

My CPU when the L1 cache misses



EECS 370

Improving Caches

Announcements

- My office hours on Wed 11/8 are cancelled
- Lab 9 meets Fr/M
- P3 Checkpoint due tonight
 - Full P3 due Thu 11/9
- Homework 3 due Mon 11/6

Agenda

- **Larger Cache Blocks**
- Extra Problems
- LRU with More than Two Blocks
- Write-Through Cache
- Write-Back Cache

Calculating Size

- How many bits is used in cache?
 - Storing data
 - 2 bytes of SRAM
 - Calculate overhead (non-data)
 - This cost is often forgotten for caches, but it drives up the cost of real designs!
 - 2 4-bit tags
 - 2 valid bits
- What is the storage requirement

Poll: Which of the following would reduce tag overhead (as an overall percentage)? (select all that apply)

- ☒ a) Increase number of cache entries *The percentage is the same.*
- ☒ b) Decrease number of cache entries
- ☒ c) Use smaller addresses *But only the ISA has the decision for the size of address.*
- ☒ d) Store more data in each cache entry

How can we reduce overhead?

- Have a smaller address
 - Impractical, and caches are supposed to be micro-architectural
- Cache bigger units than bytes
 - Each block has a single tag, and blocks can be whatever size we choose.

lru	1	1	110	
	1	7	170	
	tag		data	

Increasing Block Size

Case 1:

Block size: 1 bytes

1	0	74
1	6	160
V	tag	data (block)

How many bits needed per tag?

$$= \log_2(\text{number of blocks in memory}) = \log_2(16)$$

$$= 4 \text{ bits}$$

$$\text{Overhead} = (4+1) / 8 = 62.5\%$$

Case 2:

Block size: 2 bytes

1	0	74	110
1	3	160	170
V	tag	data (block)	

How many bits needed per tag?

$$= \log_2(\text{number of blocks in memory}) = \log_2(8)$$

$$= 3 \text{ bits}$$

$$\text{Overhead} = (3+1) / 16 = 25\%$$

advantage: ① double the data size

② decrease the tag size.

Memory	Tag (case 1)	Tag (case 2)
0	0	0
1	1	0
2	2	1
3	3	1
4	4	2
5	5	2
6	6	3
7	7	3
8	8	4
9	9	4
10	10	5
11	11	5
12	12	6
13	13	6
14	14	7
15	15	7

Figuring out the tag

- If block size is N, what's the pattern for figuring out the tag from the address?

$$\text{tag} = \left\lfloor \frac{\text{addr}}{\text{block size}} \right\rfloor$$

- If block size is power of 2, then this is just everything except the $\log_2(\text{block size})$ bits of the address in binary!

- E.g.

$$0d11 = 0b1011$$

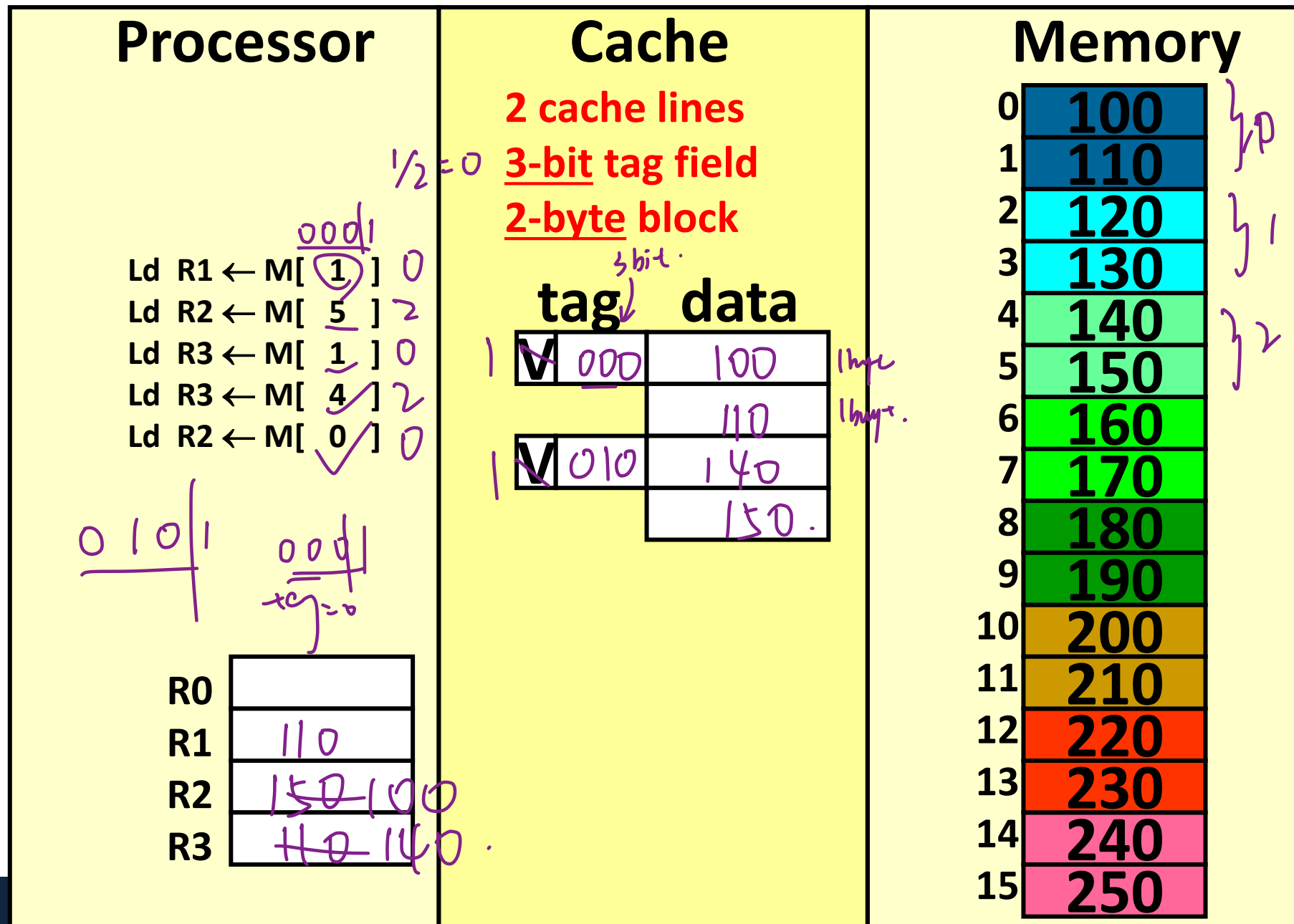
$$\text{Tag} = 0b101 = 0d5$$

$$\text{Block Offset} = 1$$

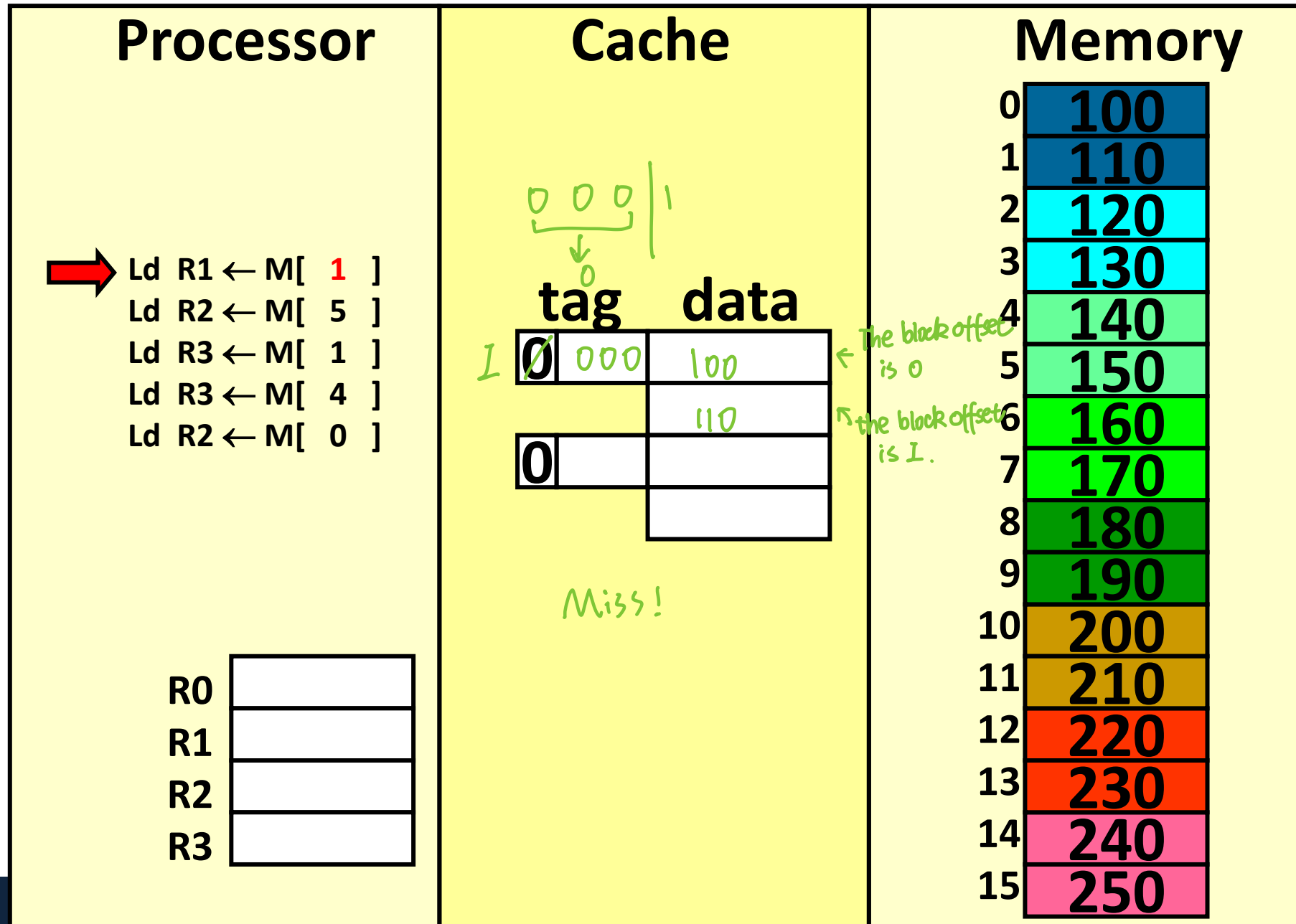
- Remaining bits (block offset) tells us how far into the block the data is

Memory		Tag (case 2)	
0	74	0	
1	110	0	
2	120	1	
3	130	1	
4	140	2	
5	150	2	
6	160	3	
7	170	3	
8	180	4	
9	190	4	
10	200	0	5
11	210	1	5
12	220	6	
13	230	6	
14	240	7	
15	250	7	

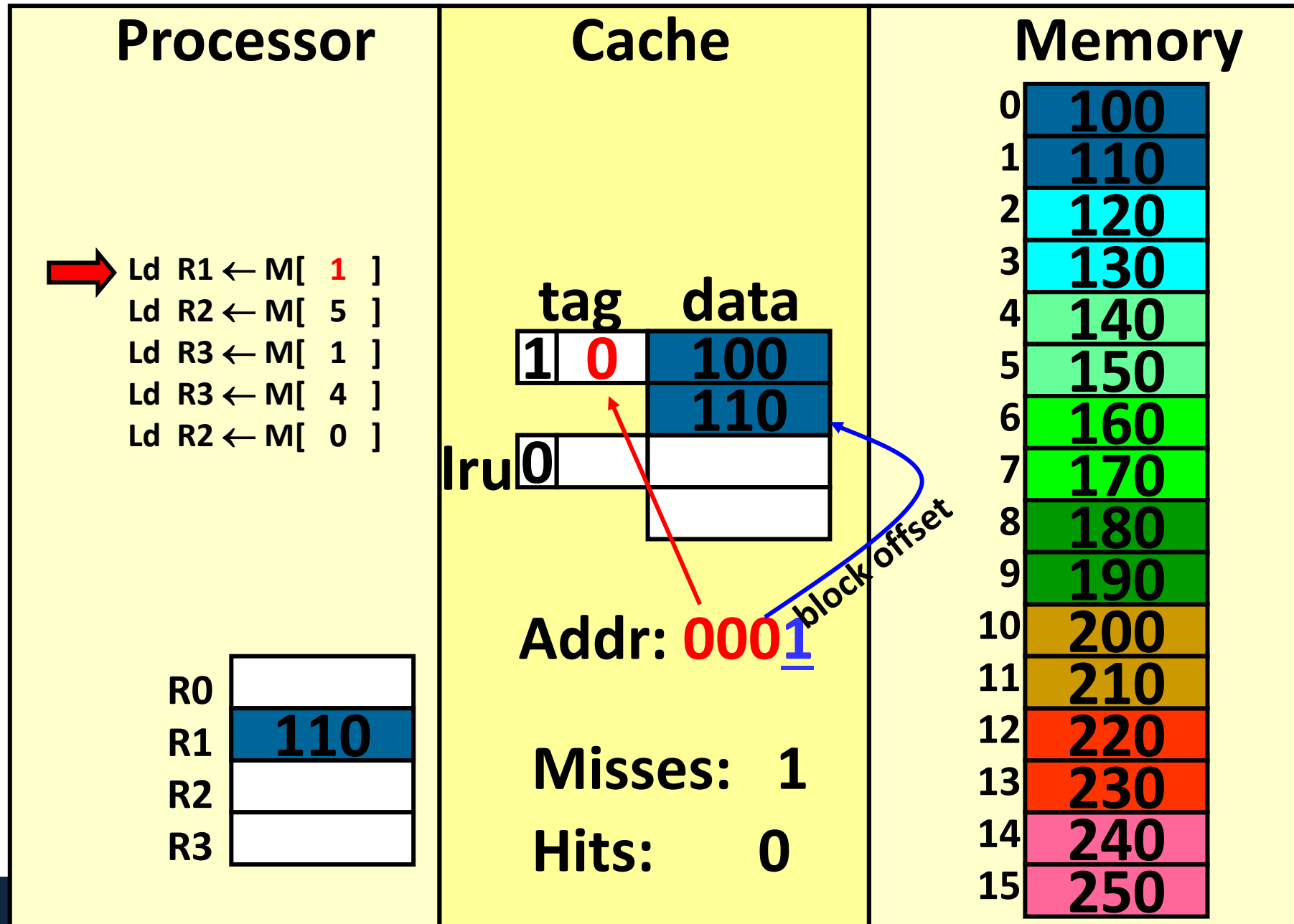
Block size for caches



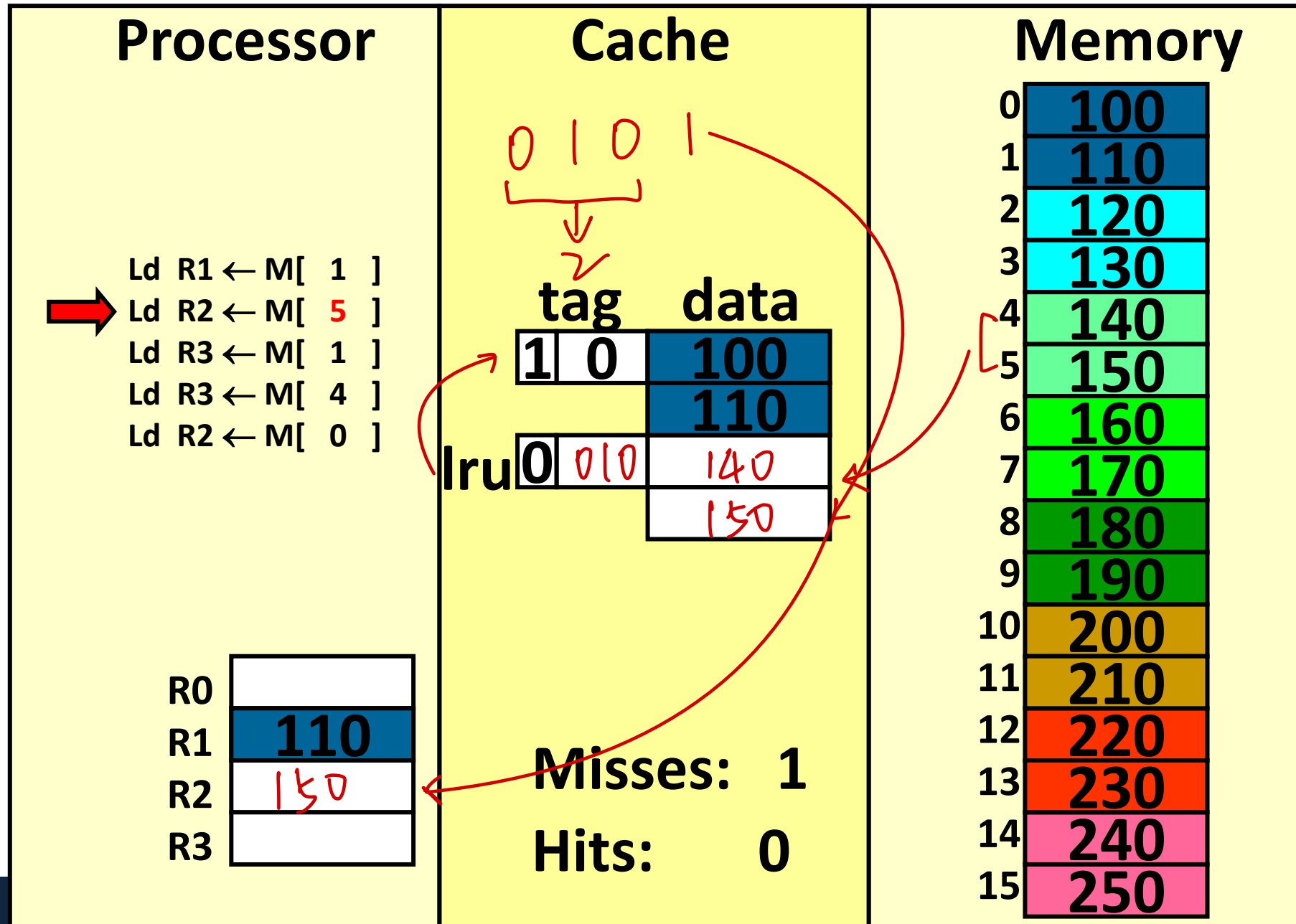
Block size for caches



Block size for caches

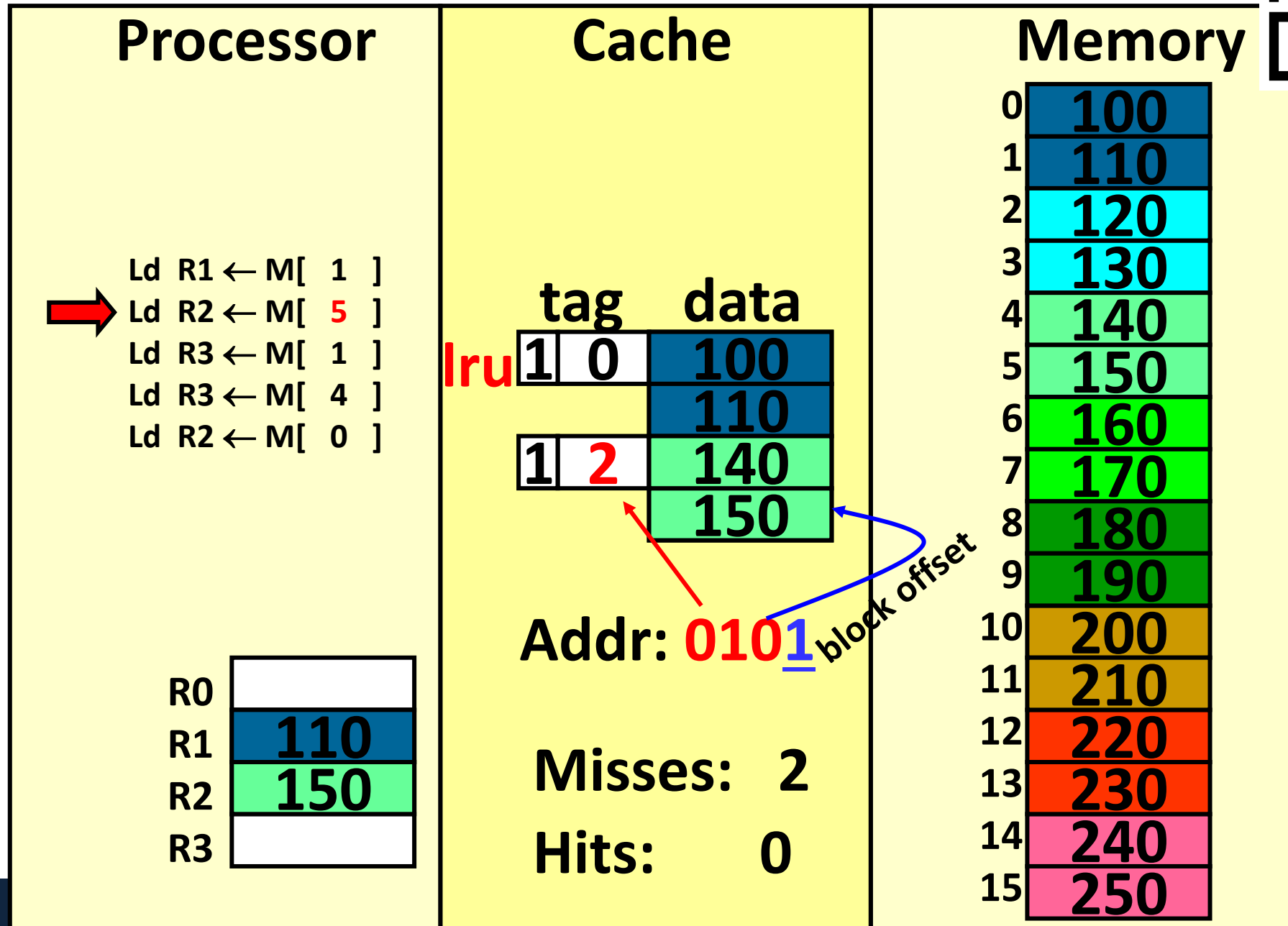
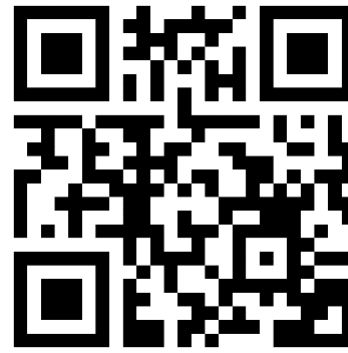


Block size for caches

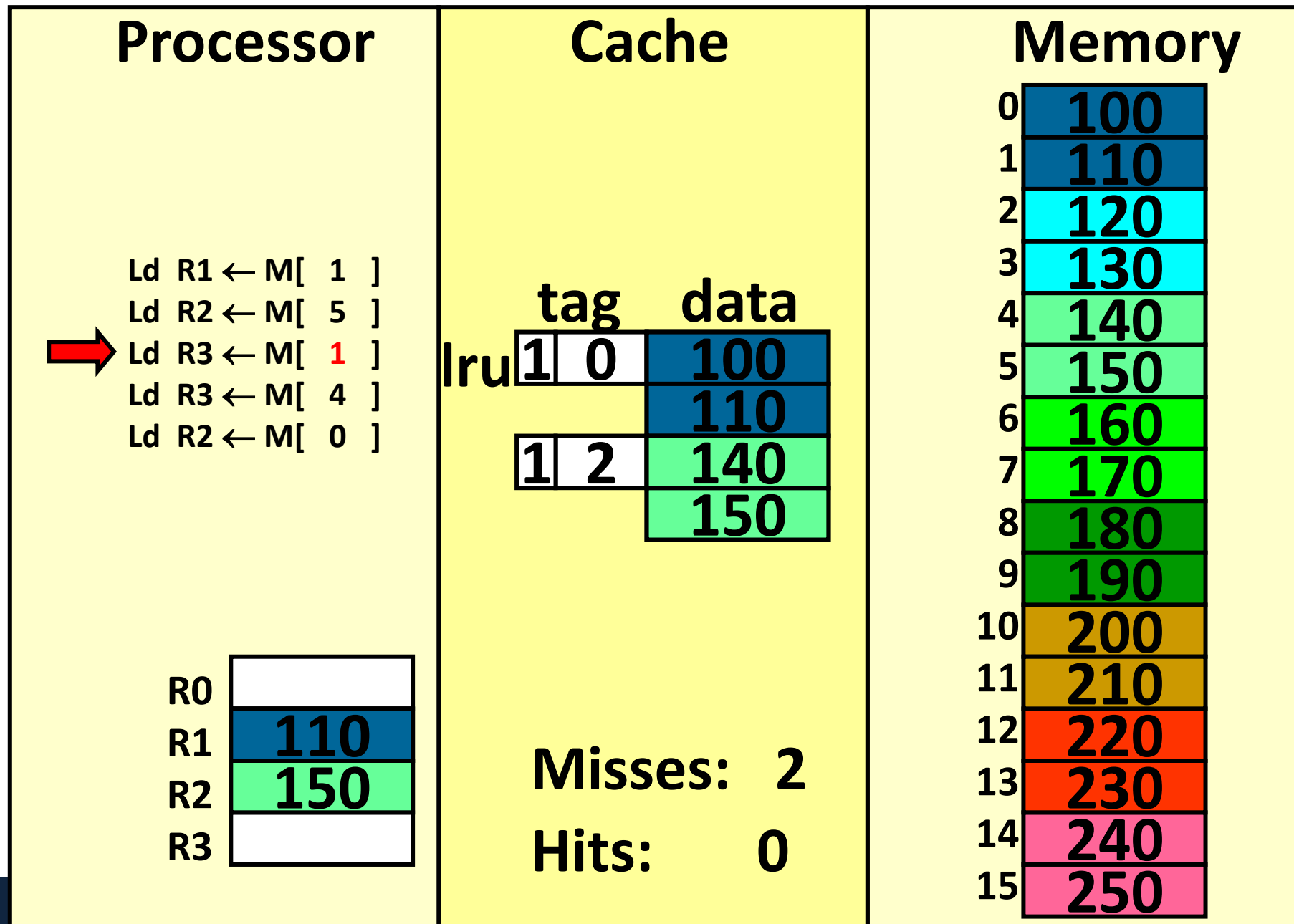


Block size for caches

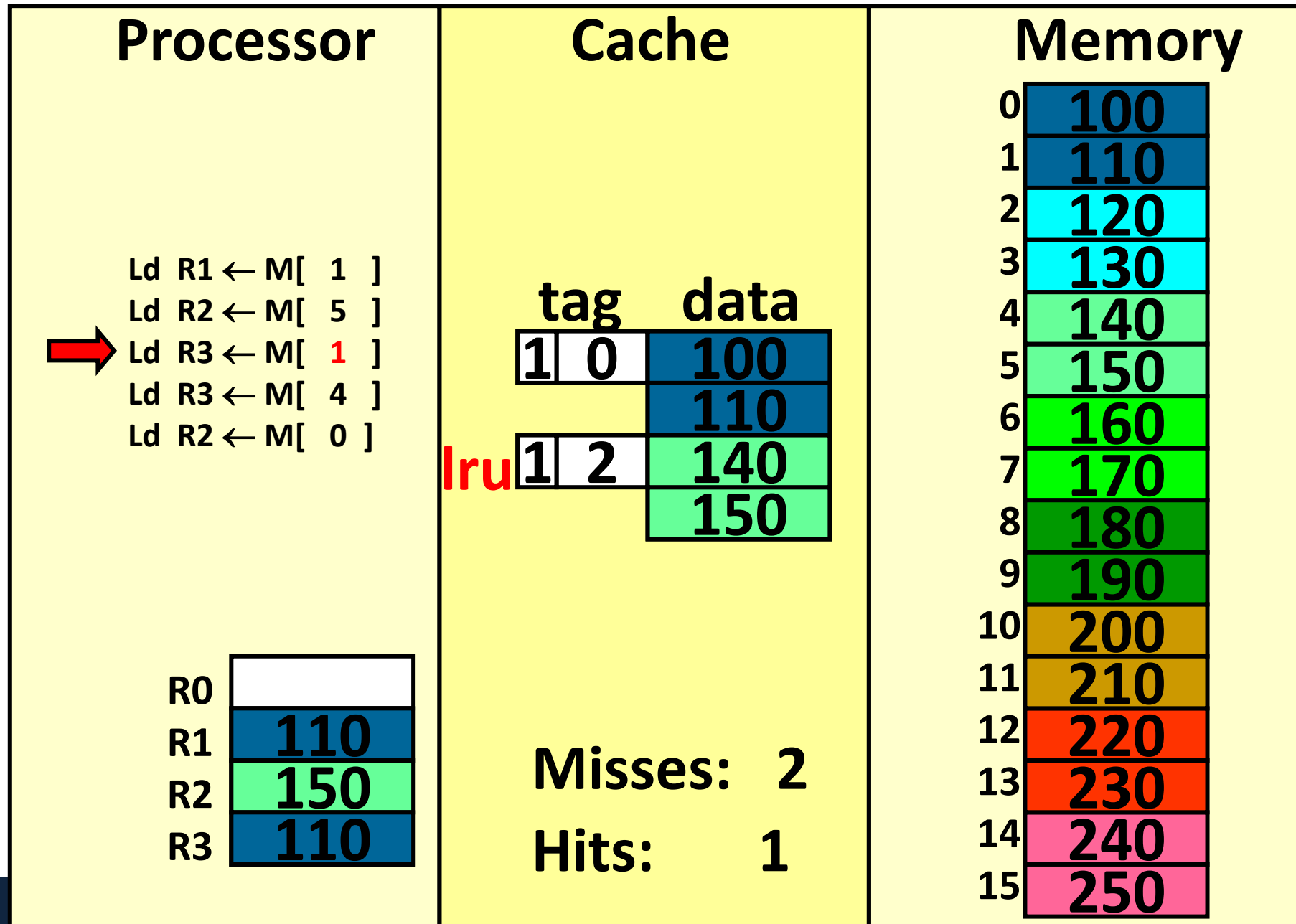
Poll: Complete the last 3 instructions yourself



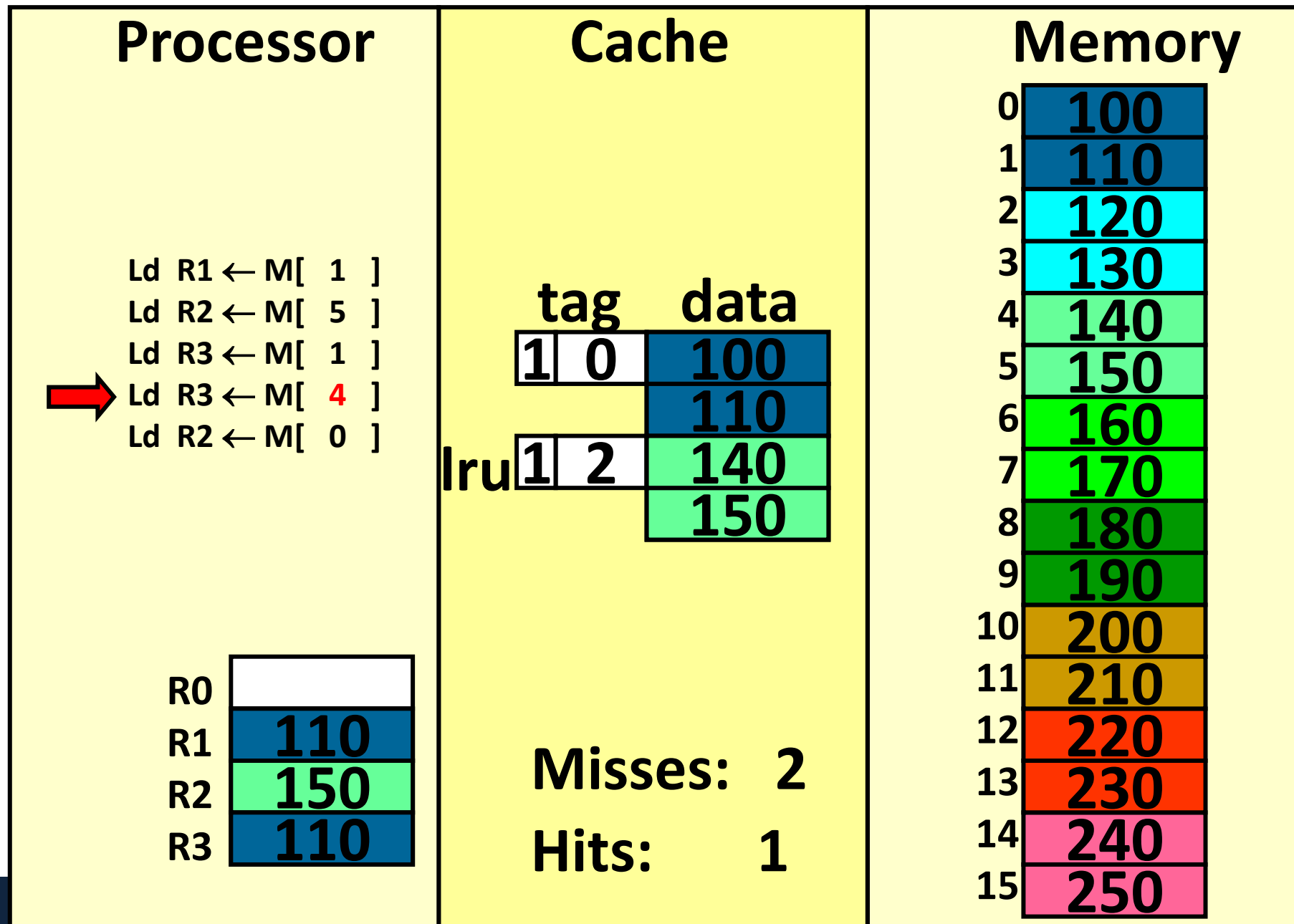
Block size for caches



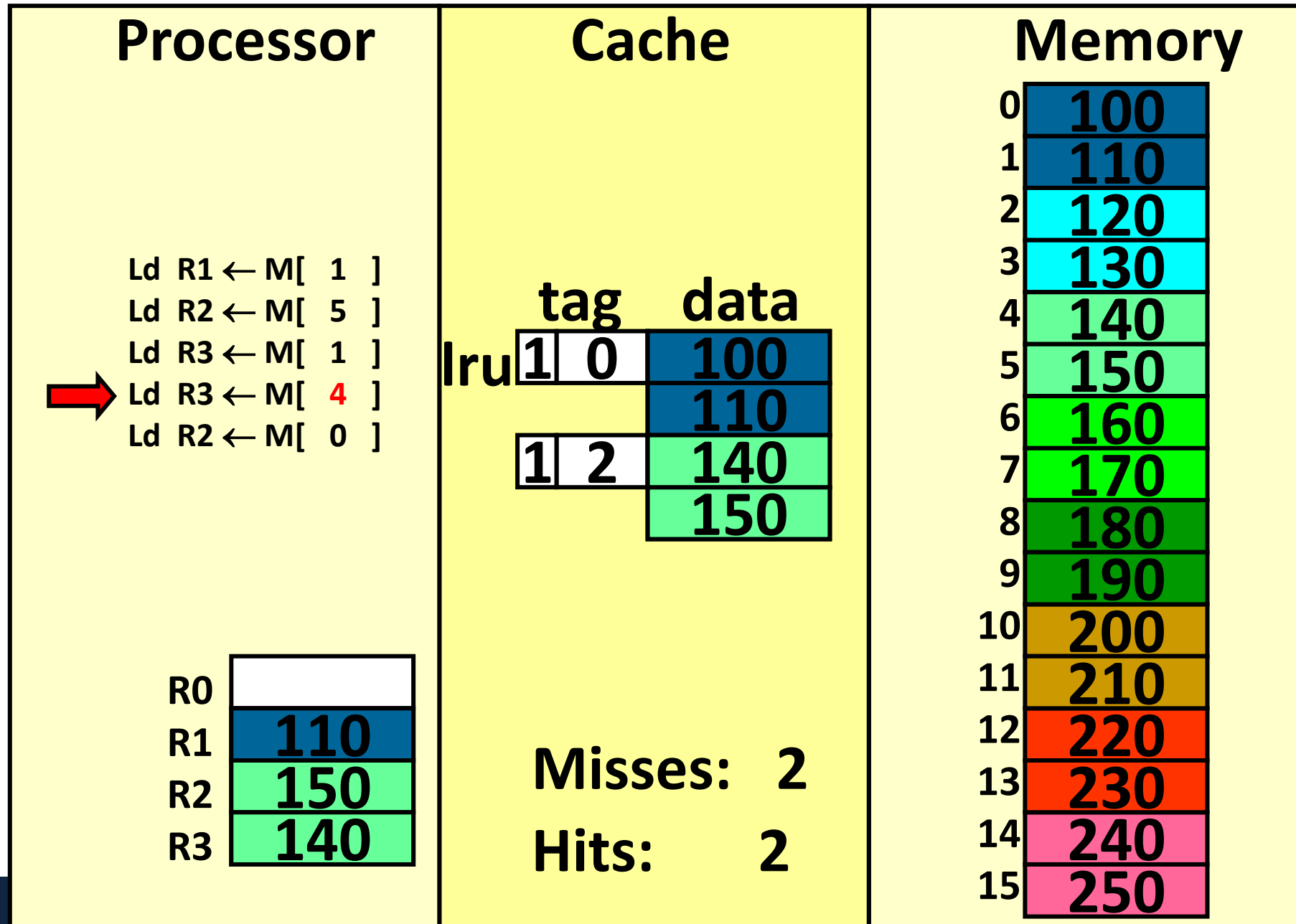
Block size for caches



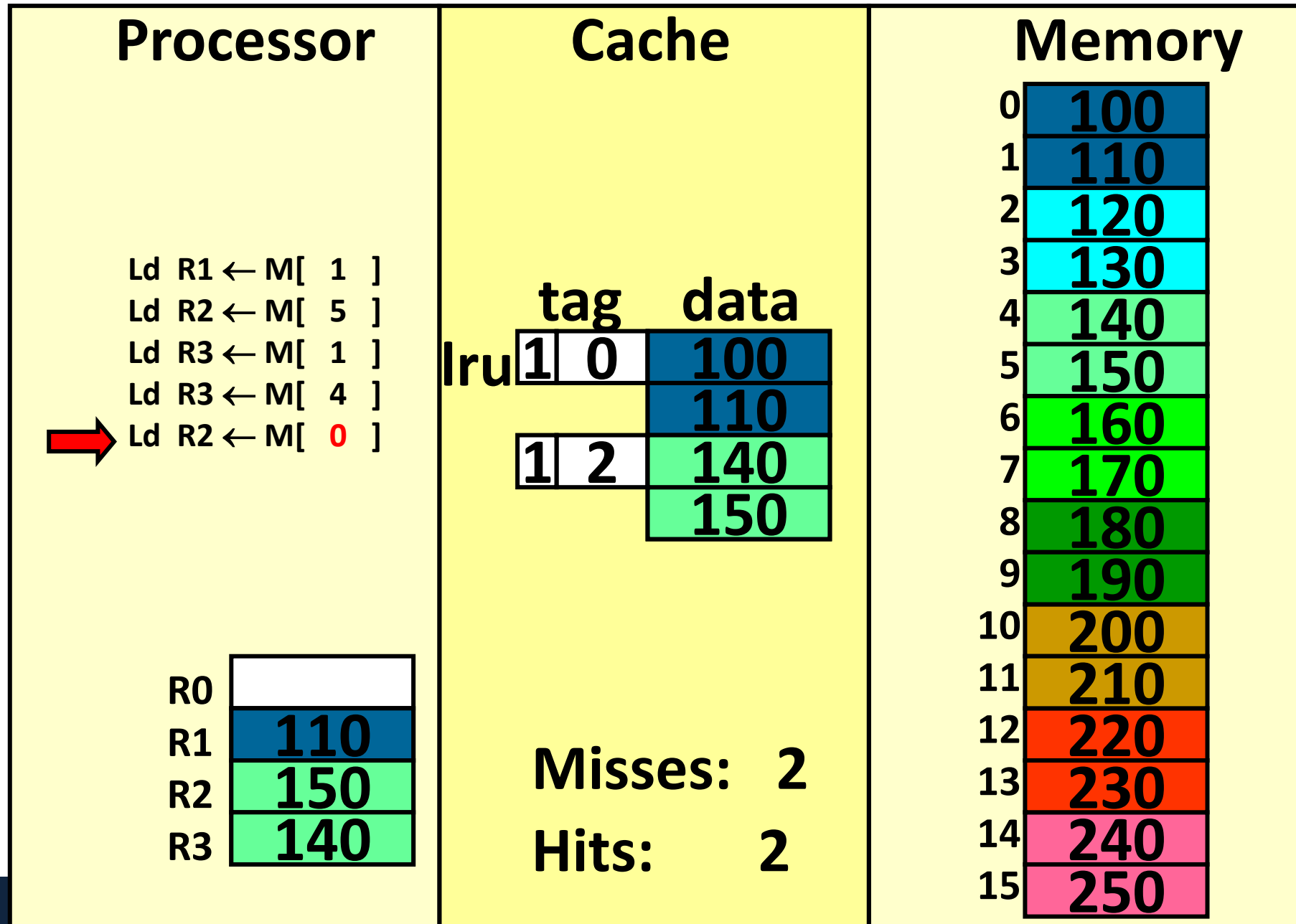
Block size for caches



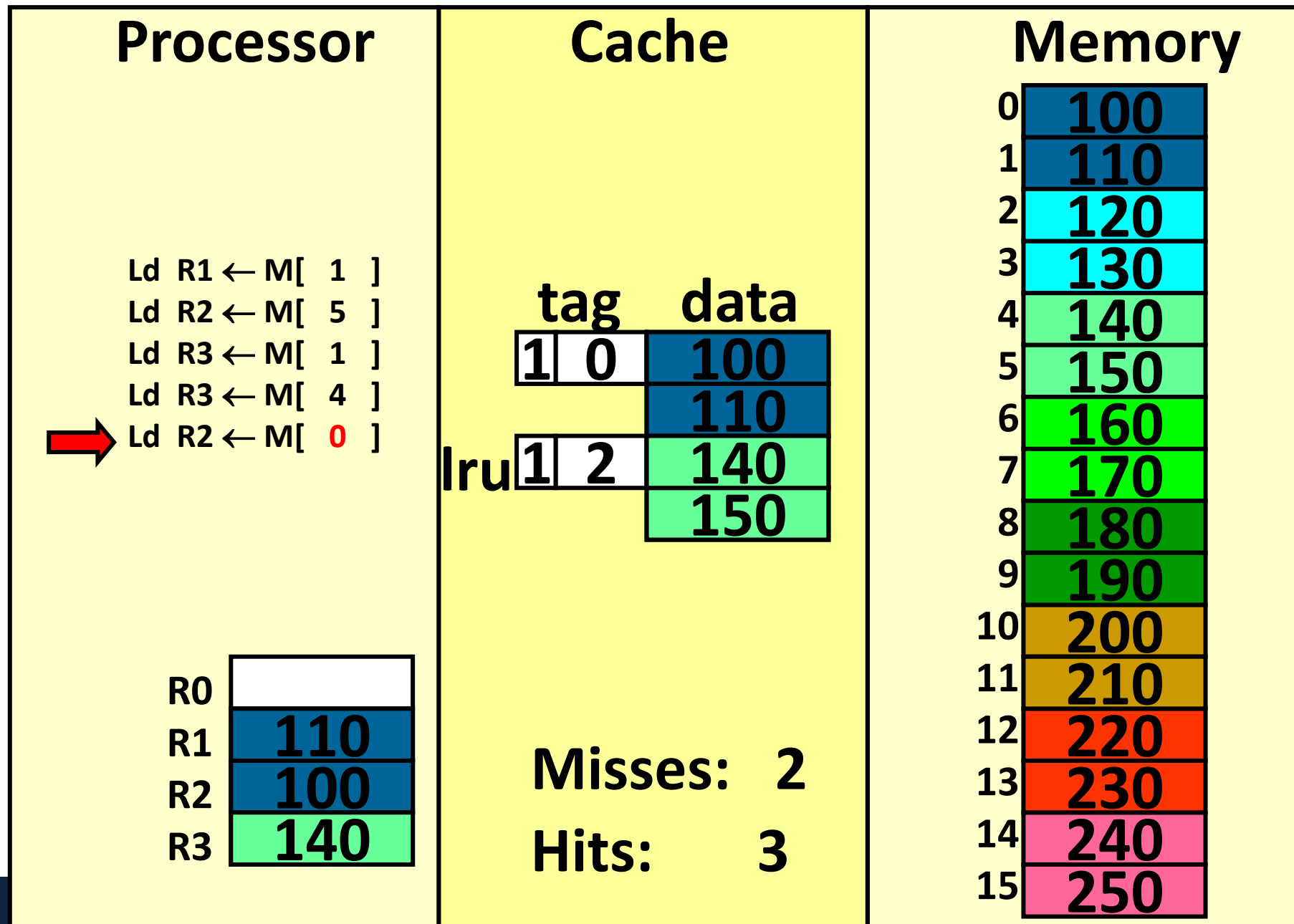
Block size for caches



Block size for caches



Block size for caches



Spatial Locality

- Notice that when we accessed address 1, we also brought in address 0
 - This turned out to be a good thing, since we later referenced address 0 and found it in the cache
- This is taking advantage of **spatial locality**:
 - If we access a memory location (e.g. 1000), we are more likely to access a location near it (e.g. 1001) than some random location
 - Arrays and structs are a big reason for this

```
for(i=0; i< N; i++)  
    for(j = 0; j < N; j++ )  
    {  
        count++;  
        arrayInt[i][j] = 10;  
    }
```

Agenda

- Larger Cache Blocks
- **Extra Problems**
- LRU with More than Two Blocks
- Write-Through Cache
- Write-Back Cache

Extra Practice Problem

Processor

0011

tagBlock offset

Ld R0 ← M[3]

Ld R2 ← M[12]

Ld R3 ← M[15]

Ld R1 ← M[4]

Ld R2 ← M[9]

0011

R0

R1

R2

R3

123

71

221

225

Cache

2 cache lines

2-bit tag field

4-byte block

Vtag data

1001

1032

7871

29150

172162

173173

1818

2121

2833

22528

Memory

078

129

2120

3123

471

5150

6162

7173

818

921

1033

1128

1219

13200

14210

15225

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

0

1

2

3



Solution to Practice Problem

Ld R0 \leftarrow M[3]

Ld R2 \leftarrow M[12]

Ld R3 \leftarrow M[15]

Ld R1 \leftarrow M[4]

Ld R2 \leftarrow M[9]

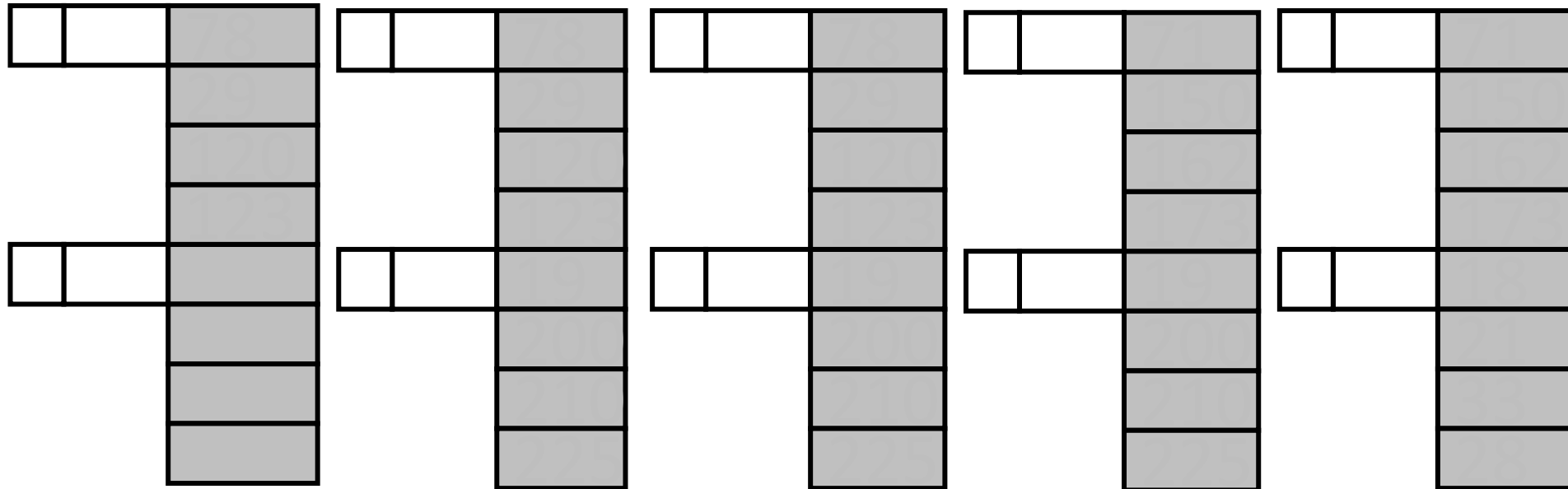
V tag data

V tag data

V tag data

V tag data

V tag data



Solution to Practice Problem

Ld R0 ← M[3]

Ld R2 ← M[12]

Ld R3 ← M[15]

Ld R1 ← M[4]

Ld R2 ← M[9]

V tag data

V tag data

V tag data

V tag data

V tag data

1	0	78
		29
		120
		123
0		

lru

miss

1	0	78
		29
		120
		123
1	3	19
		200
		210
		225

lru

miss

1	0	78
		29
		120
		123
1	3	19
		200
		210
		225

lru

hit

1	1	71
		150
		162
		173
1	3	19
		200
		210
		225

lru

miss

1	1	71
		150
		162
		173
1	2	18
		21
		33
		28

lru

miss

Extra Class Problem

*We'll see later that this is called a "fully-associative cache"

- Given a cache that works as we've described* with the following configuration: total size is 8 bytes, block size is 2 bytes, LRU replacement. The memory address size is 16 bits and is byte addressable.

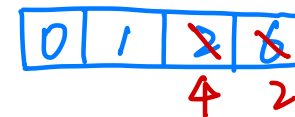
- How many bits are for each tag? How many blocks in the cache?

Tag = 16 - (block offset) = 15 bits;

2 byte block, 8 bytes total = 4 blocks.

- For the following reference stream, indicate whether each reference is a hit or miss: 0, 1, 3, 5, 12, 1, 2, 9, 4

1 2 6 0 1 4 2
M M M H H M M



0 0
M H

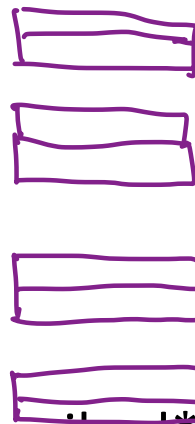
- What is the hit rate?

$3/9 = 33\%$

- How many bits are needed for storage overhead for each block?

overhead = Tag + Valid + LRU = 18 bits
(15) (1) ($\log_2^{(4)}$)

Extra Class Problem



2^{16} memory.

*We'll see later that this is called a "fully-associative cache"

- Given a cache that works as we've described* with the following configuration: total size is 8 bytes, block size is 2 bytes, LRU replacement. The memory address size is 16 bits and is byte addressable.

1. How many bits are for each tag? How many blocks in the cache?

$$16 - 1 = 15$$

15

4.

2. For the following reference stream, indicate whether each reference is a hit or miss: 0, 1,

3, 5, 12, 1, 2, 9, 4

m m m h h m

0
1
2
3
4 8
5 9
6 12
7 1
8 2
9 9
10 4

miss. hit.

3. What is the hit rate?

miss: 6
hit: 3
 $3/9 = \frac{1}{3}$

4. How many bits are needed for storage overhead for each block?

$$\frac{15}{\text{tag}} + \frac{\text{LRU}}{2} + \frac{\text{Validation}}{1}$$

18

Agenda

- Larger Cache Blocks
- Extra Problems
- **LRU with More than Two Blocks**
- Write-Through Cache
- Write-Back Cache

LRU with more than 2 entries

- If we have more than 2 things we're keeping track of...
 - Can't **just** track LRU
 - Once we access that element, how do we know which of the other elements are LRU?
 - Must track the *full ordering* of when elements were accessed*
- Each element must store a number $[0-(N-1)] \rightarrow \log_2(N)$ bits
- 0 is LRU, 1 is 2nd LRU... N-1 is most recently used

LRU with more than 2 entries

- If we have more than 2 things we're keeping track of...
 - Can't **just** track LRU
 - Once we access that element, how do we know which of the other elements are LRU?
 - Must track the *full ordering* of when elements were accessed*
- Each element must store a number $[0-(N-1)] \rightarrow \log_2(N)$ bits
- 0 is LRU, 1 is 2nd LRU... N-1 is most recently used
- When element i is used:
 $X = \text{counter}[i]$
 $\text{counter}[i] = N-1$
 for ($j=0$ to $N-1$)
 if ($(j \neq i) \text{ AND } (\text{counter}[j] > X)$) $\text{counter}[j] -$

Initial State				
Element	0	1	2	3
<u>Count</u>	0	1	2	3

least recently use

most recently use.

LRU with more than 2 entries

- If we have more than 2 things we're keeping track of...
 - Can't **just** track LRU
 - Once we access that element, how do we know which of the other elements are LRU?
 - Must track the *full ordering* of when elements were accessed*
- Each element must store a number $[0-(N-1)] \rightarrow \log_2(N)$ bits
- 0 is LRU, 1 is 2nd LRU... N-1 is most recently used
- When element i is used:
 - X = counter[i]
 - counter[i] = N-1
 - for (j=0 to N-1)
 - if ((j != i) AND (counter[j] > X)) counter[j]—

Initial State				
Element	0	1	2	3
Count	0	1	2	3

Access Element 2				
Element	0	1	2	3
Count	0	1	3	2

Find the Count for all Element, and check if it's bigger than 2.
if it is, then "-1"

LRU with more than 2 entries

- If we have more than 2 things we're keeping track of...
 - Can't **just** track LRU
 - Once we access that element, how do we know which of the other elements are LRU?
 - Must track the *full ordering* of when elements were accessed*
- Each element must store a number $[0-(N-1)] \rightarrow \log_2(N)$ bits
- 0 is LRU, 1 is 2nd LRU... N-1 is most recently used
- When element i is used:
 - $X = \text{counter}[i]$
 - $\text{counter}[i] = N-1$
 - for ($j=0$ to $N-1$)
 - if ($(j \neq i) \text{ AND } (\text{counter}[j] > X)$) $\text{counter}[j] -$
- Evict element with counter = 0 when needed
- Get's expensive for moderate to large N

Initial State				
Element	0	1	2	3
Count	0	1	2	3

Access Element 2				
Element	0	1	2	3
Count	0	<u>1</u>	<u>3</u>	<u>2</u>

Access Element 0				
Element	0	1	2	3
Count	<u>3</u>	0	2	1

Agenda

- Larger Cache Blocks
- Extra Problems
- LRU with More than Two Blocks
- **Write-Through Cache**
- Write-Back Cache

What about stores?

- Where should you write the result of a store?

- If that memory location is in the cache:

- choice 1 {
- Send it to the cache.
 - Should we also send it to memory?

(write-through policy)

- If it is not in the cache:

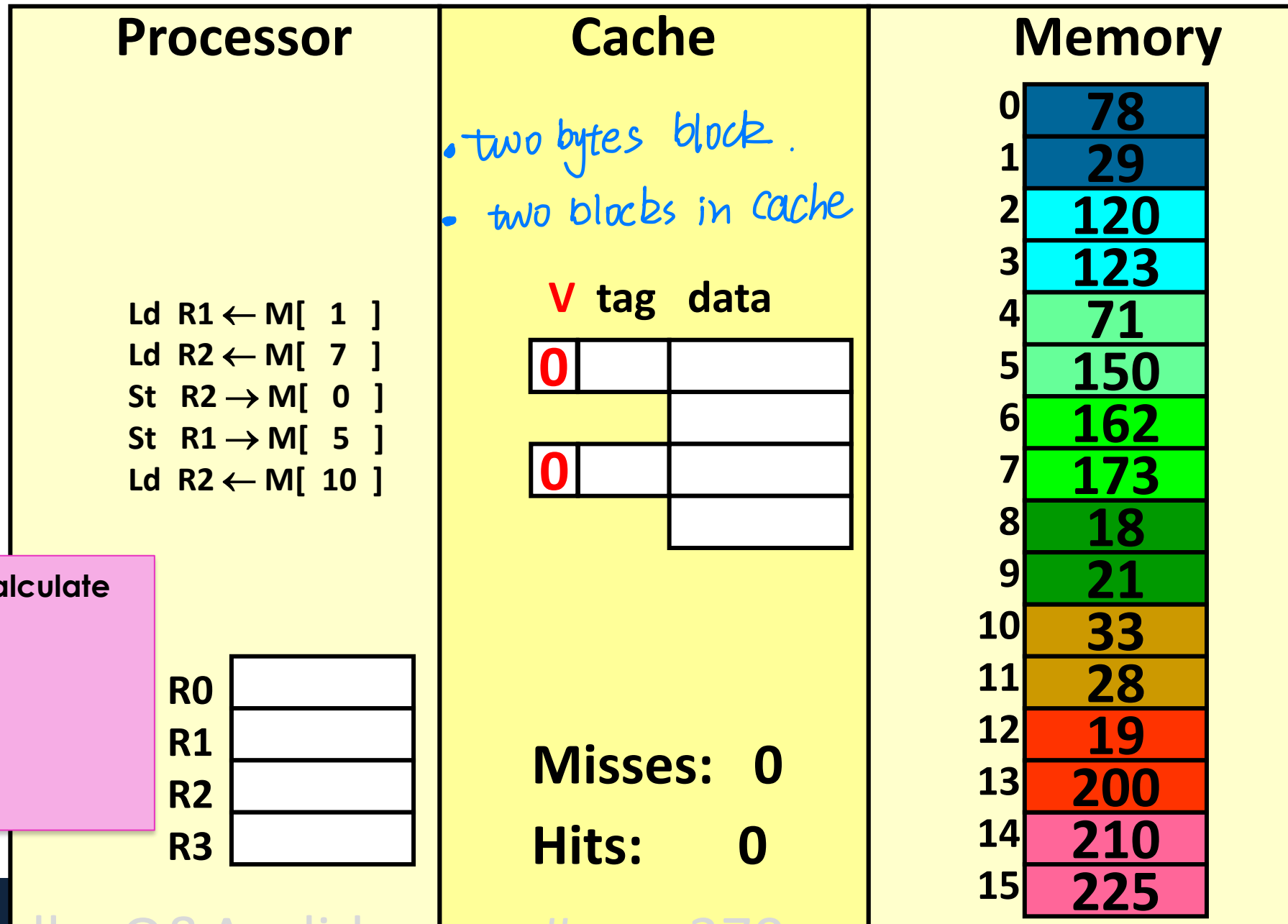
- choice 2 {
- Allocate the line (put it in the cache)?
 - Write it directly to memory without allocation?

(allocate-on-write policy)

(no allocate-on-write policy)

choice 1 and choice 2 are independent to each other

Handling stores (write-through, allocate on write)

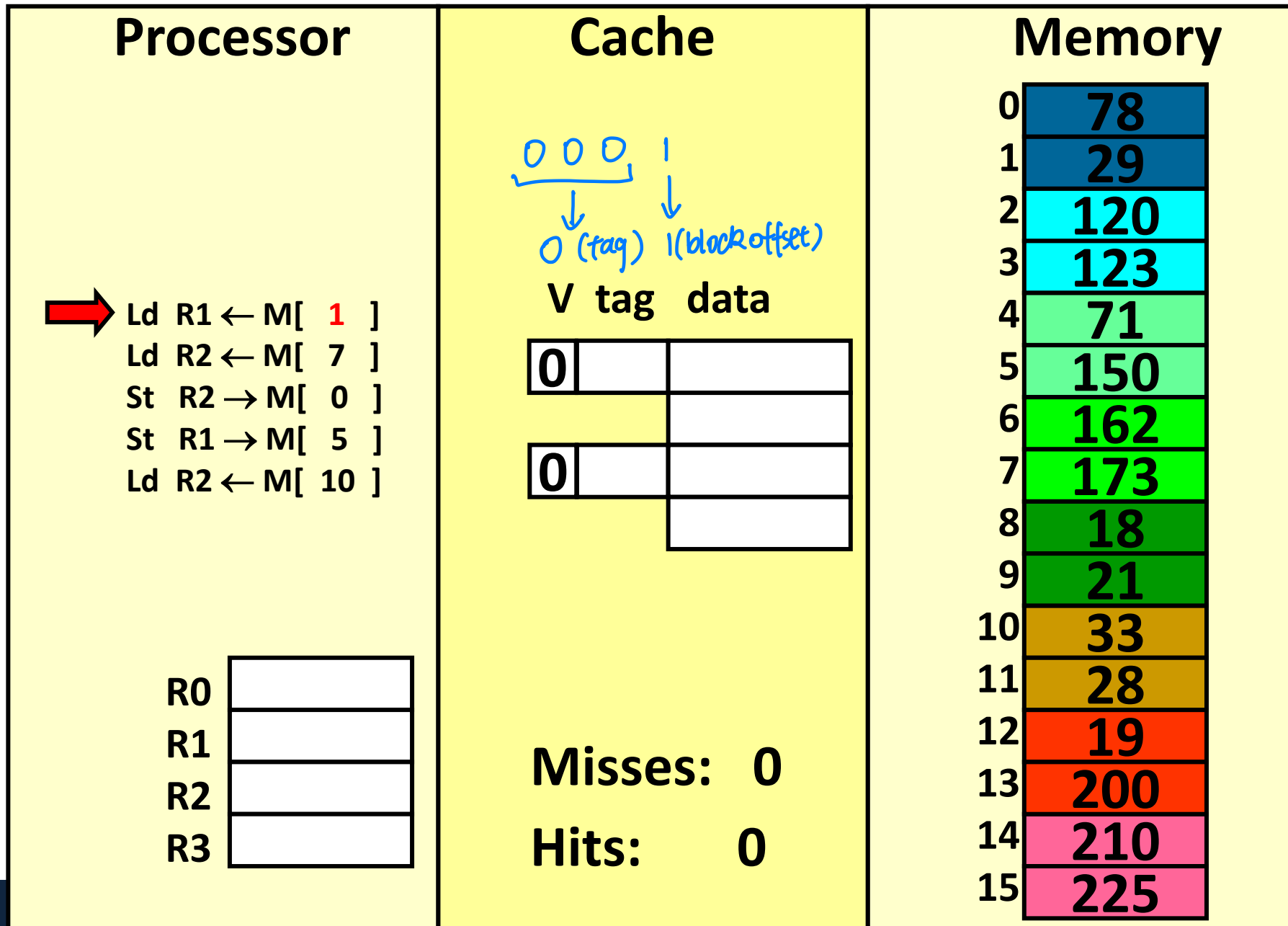


Poll: How can we calculate the tag? *b*

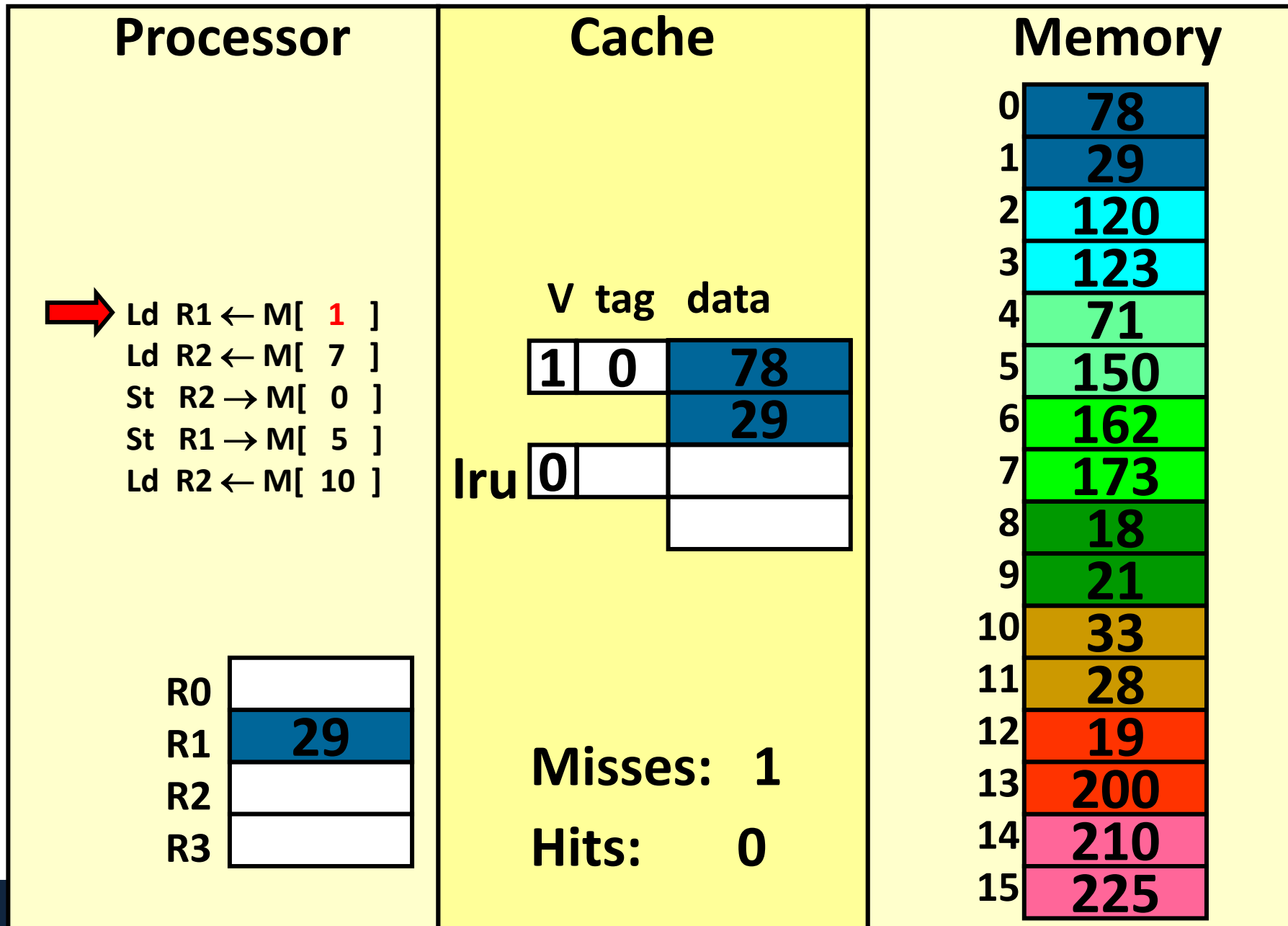
- a) addr
- b) $\text{floor}(\text{addr}/2)$ ✓
- c) $\text{ceil}(\text{addr}/2)$
- d) $\text{addr} * 2$

R0	
R1	
R2	
R3	

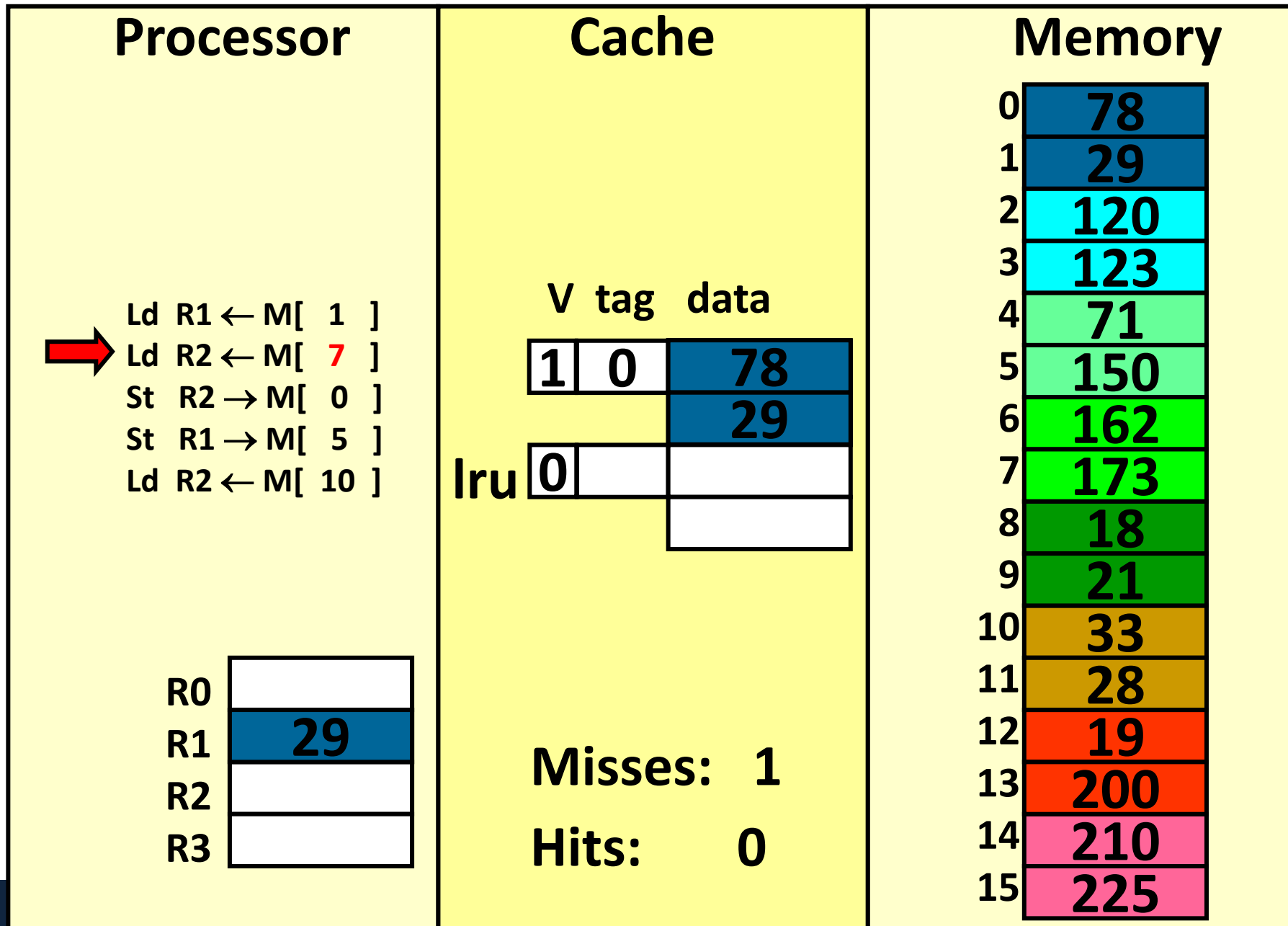
write-through, allocate on write (REF 1)



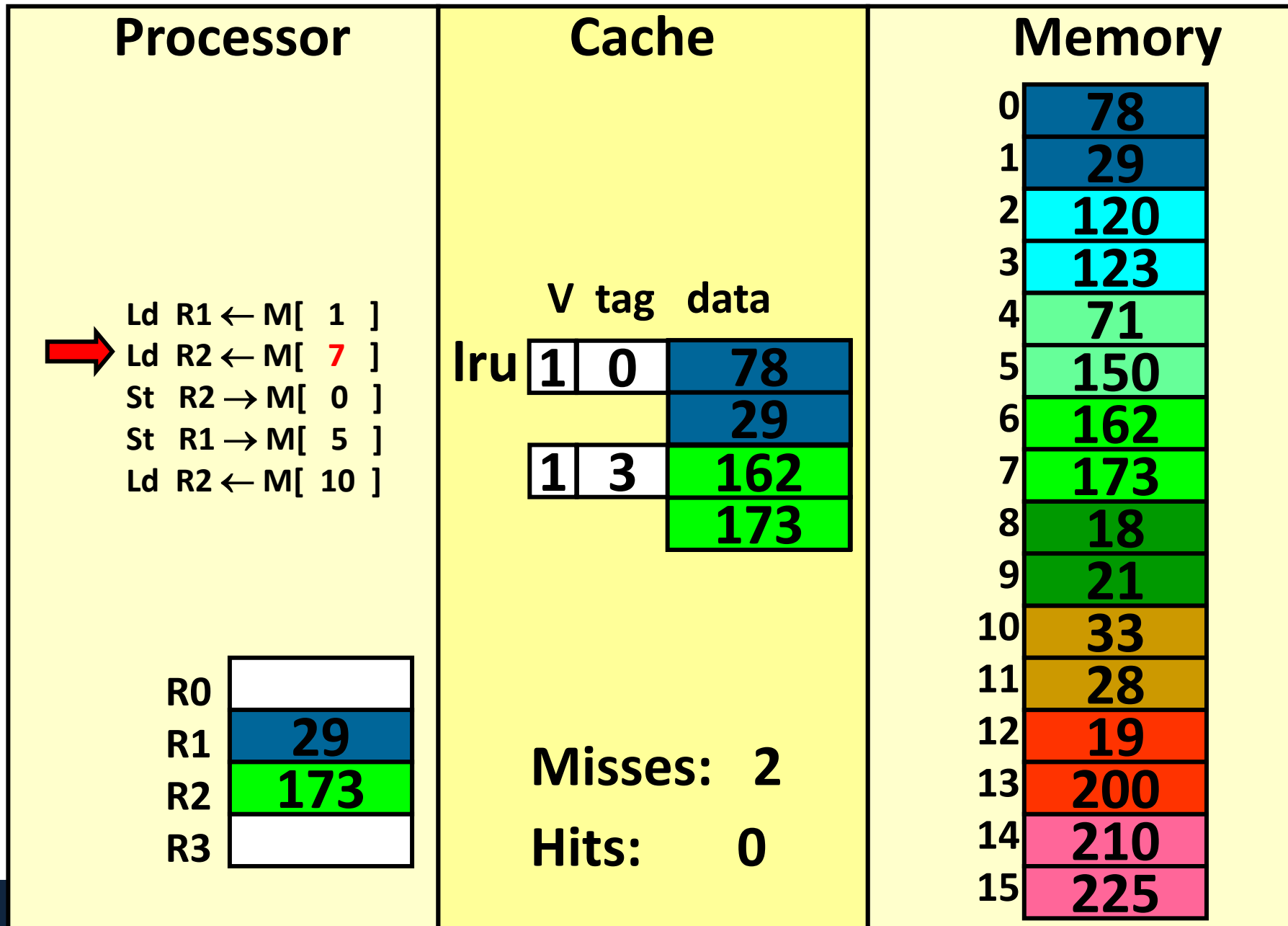
write-through, allocate on write (REF 1)



write-through, allocate on write (REF 2)

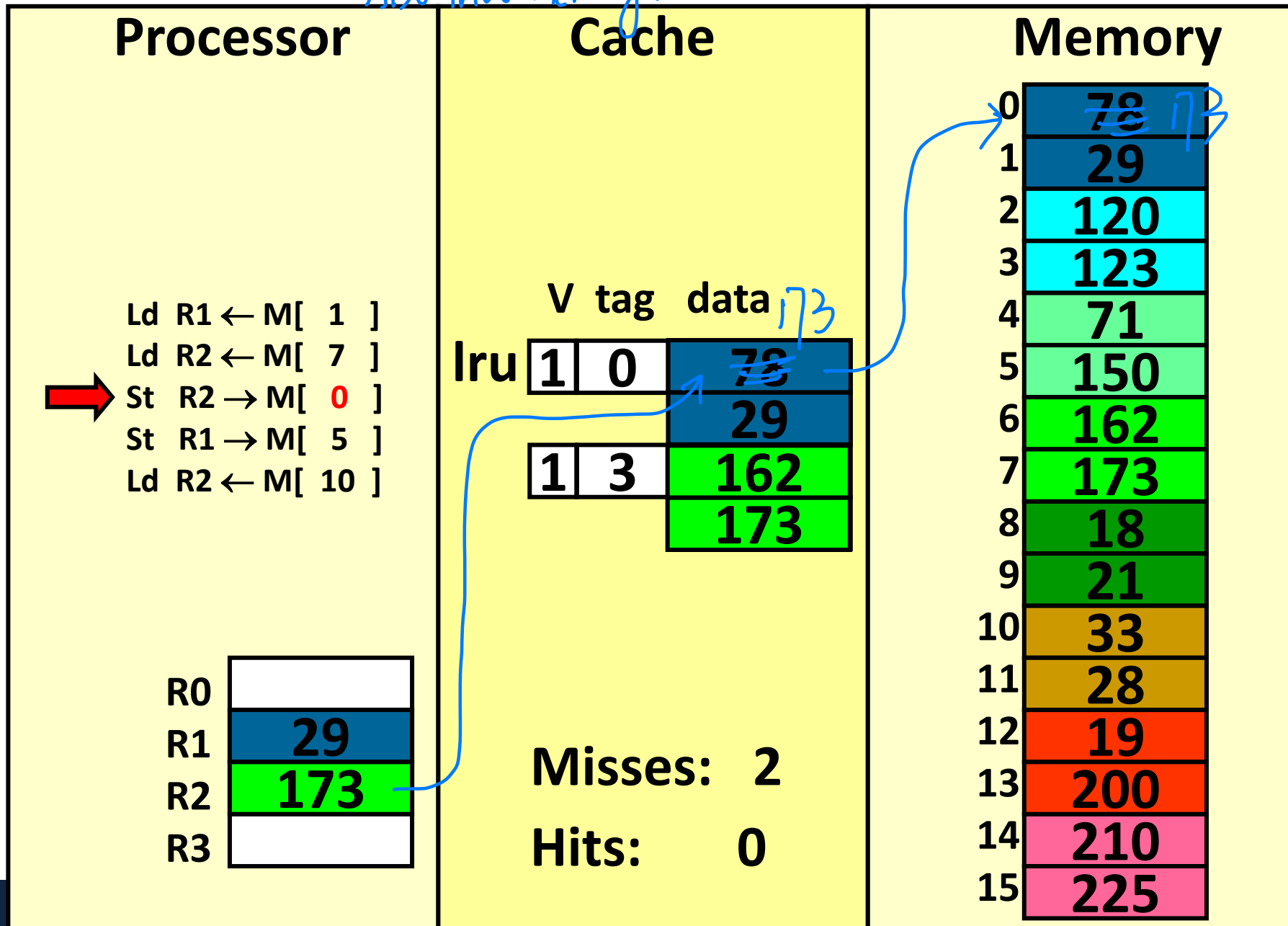


write-through, allocate on write (REF 2)

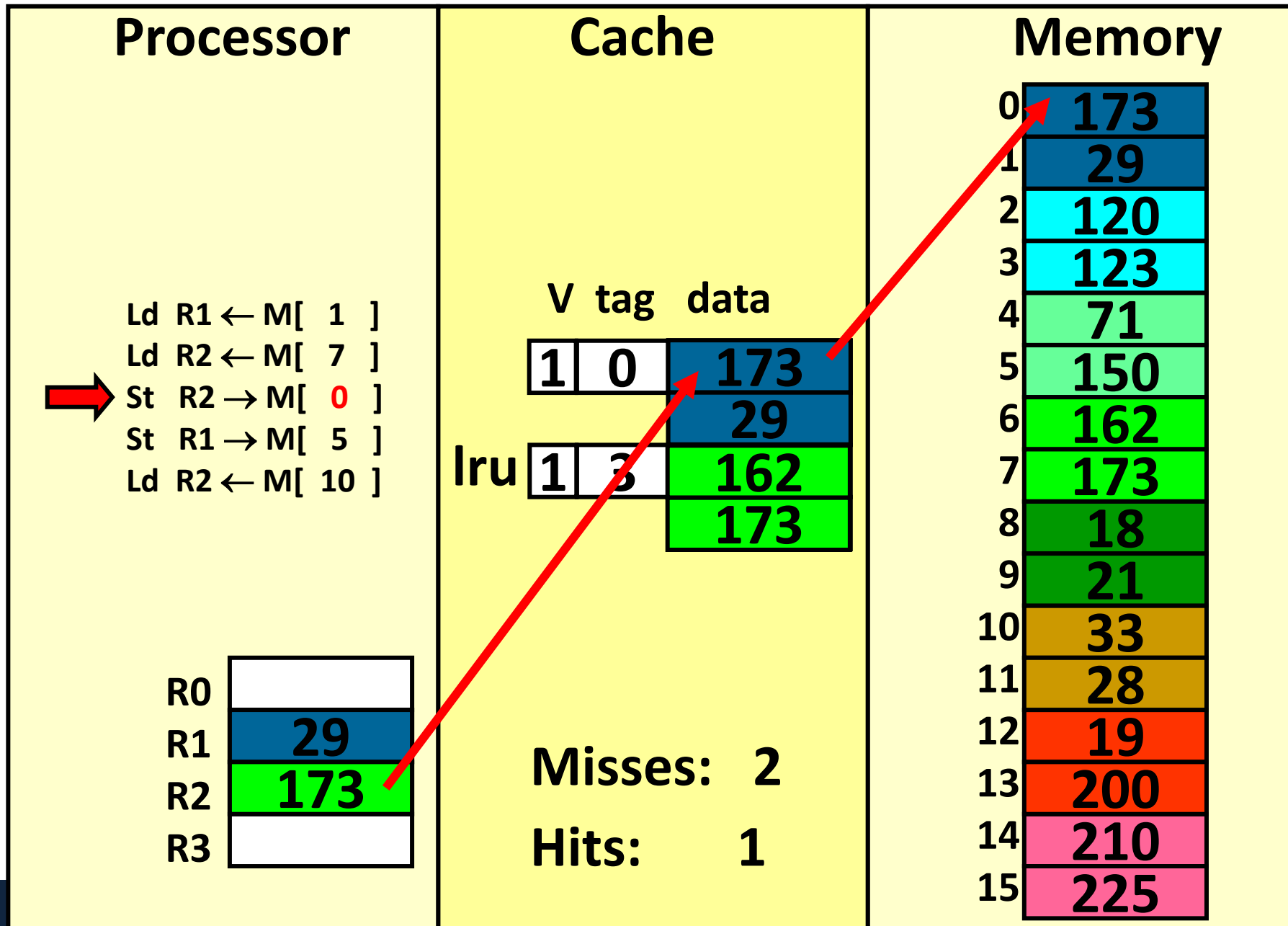


write-through, allocate on write (REF 3)

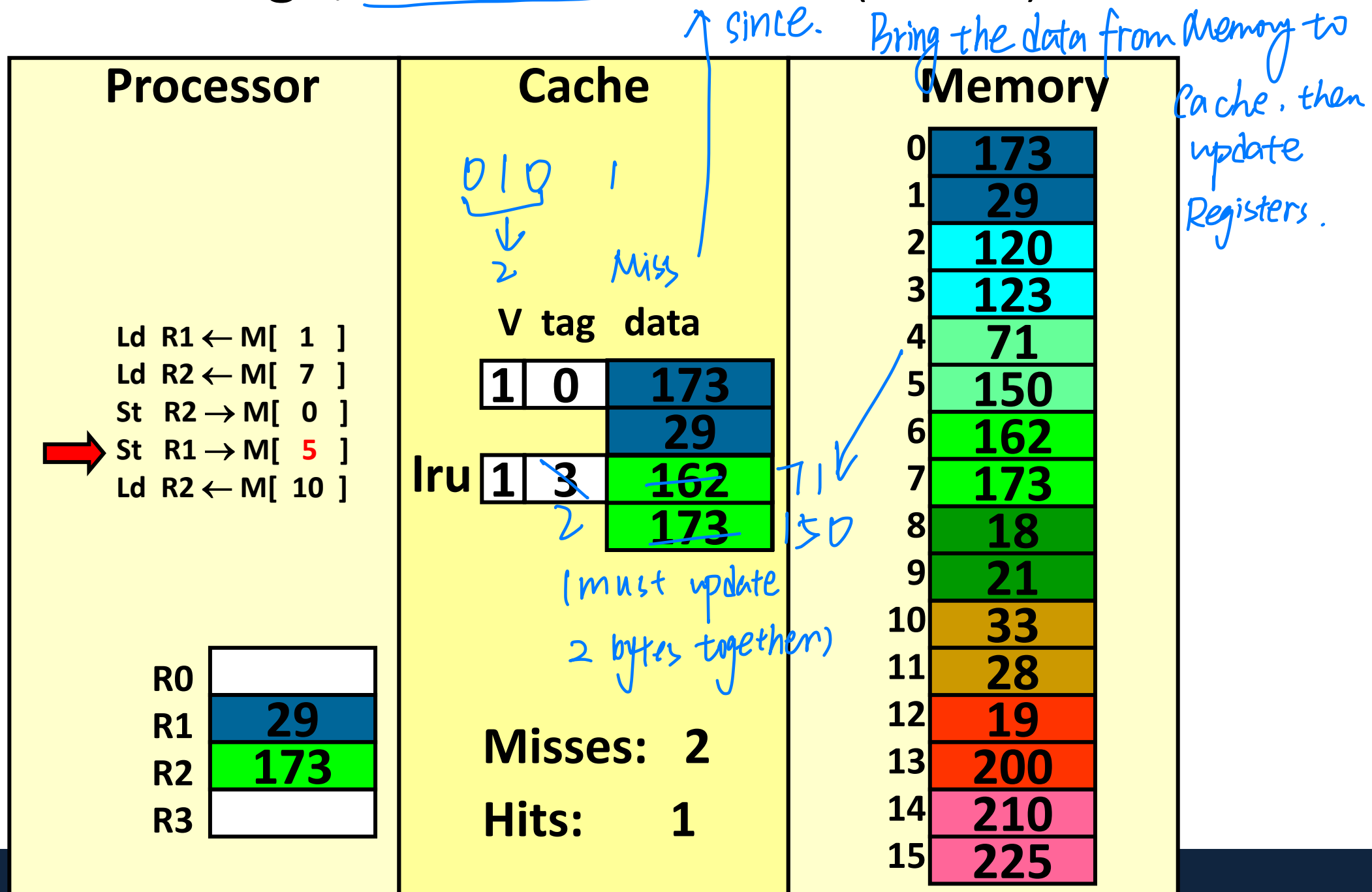
Also into memory.



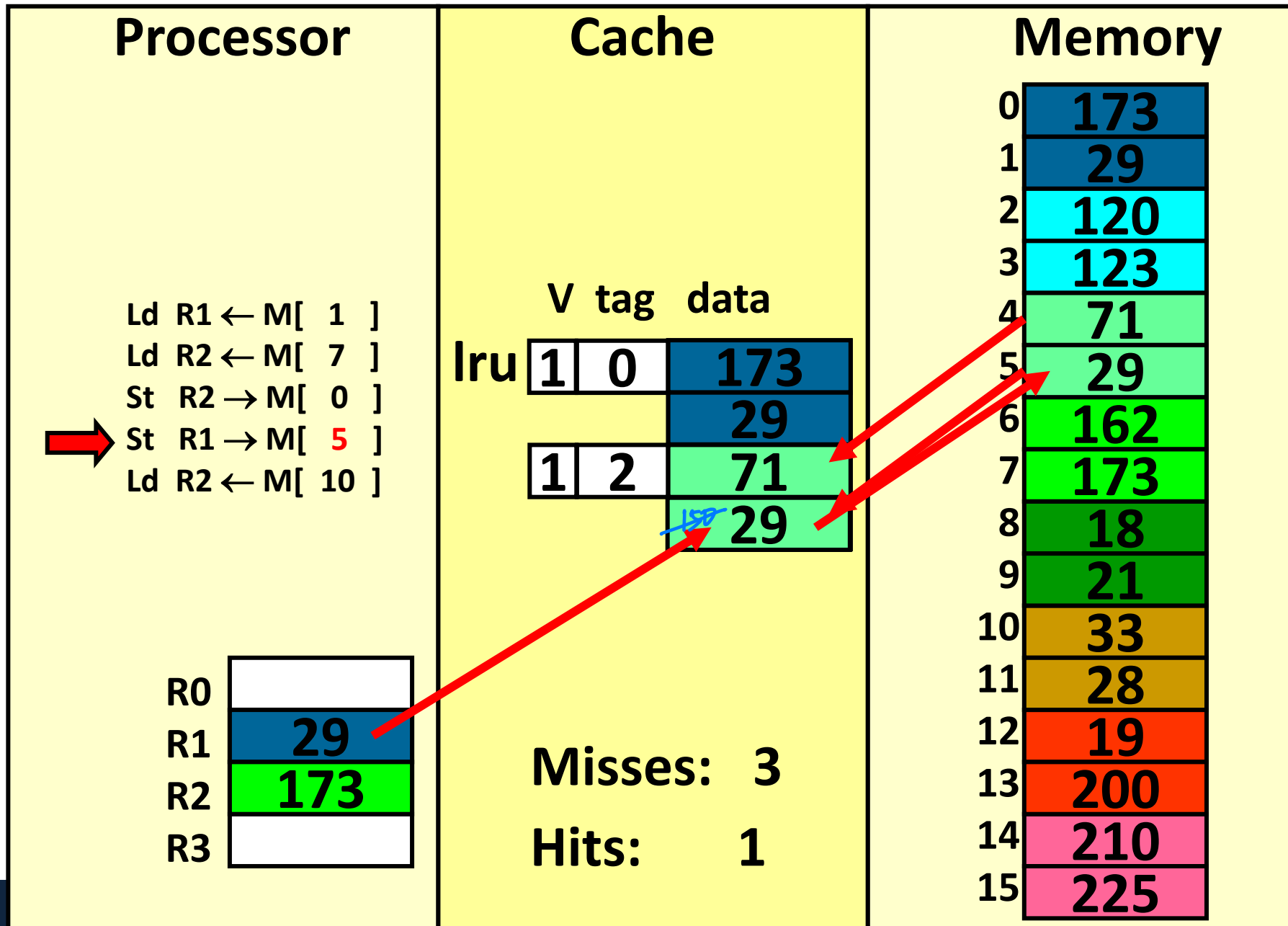
write-through, allocate on write (REF 3)



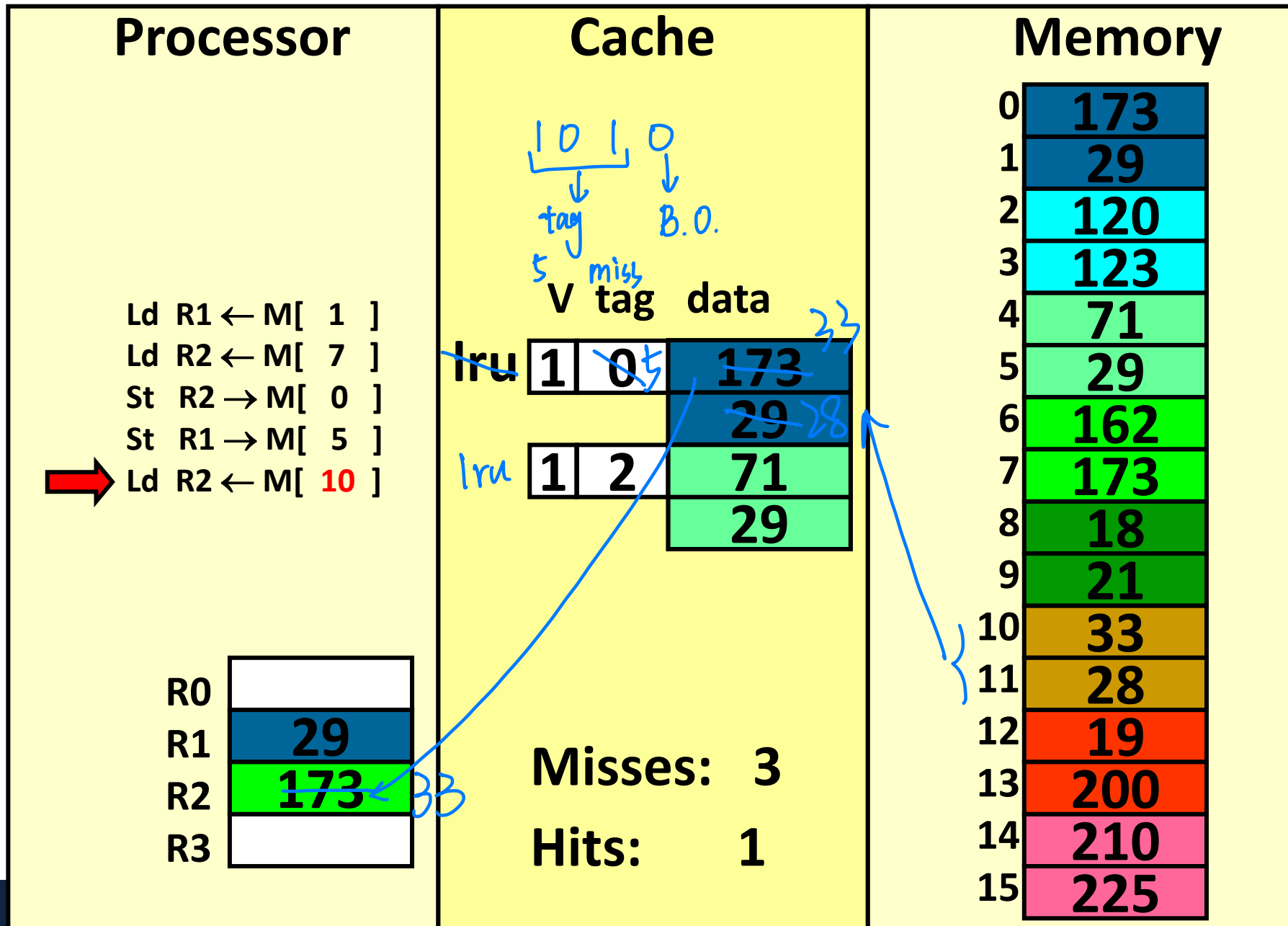
write-through, allocate on write (REF 4)



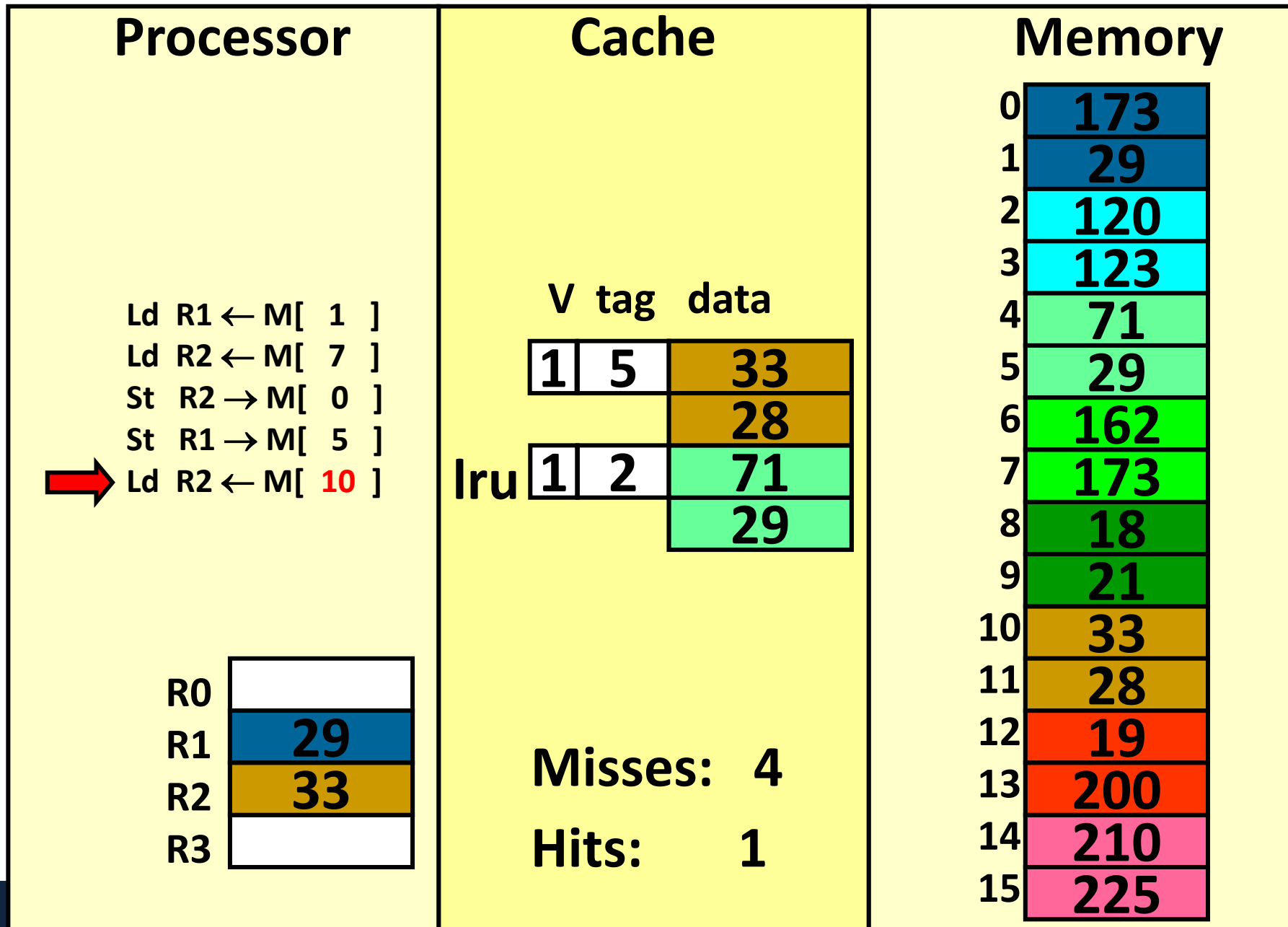
write-through, allocate on write (REF 4)



write-through, allocate on write (REF 6)



write-through, allocate on write (REF 6)



How many memory references?

- Each miss reads a block
 - 2 bytes in this cache
- Each store writes a byte
- Total reads: 8 bytes *Miss = 4 4x2 = 8 .*
- Total writes: 2 bytes *2 store instruction 2x1 = 2*
- but caches generally miss < 20%
 - Can we take advantage of that?
 - Multi-core processors have limited bandwidth between caches and memory
 - Extra stores also cost power

Next time

- Write-back Caches
- Direct-mapped vs associative caches.

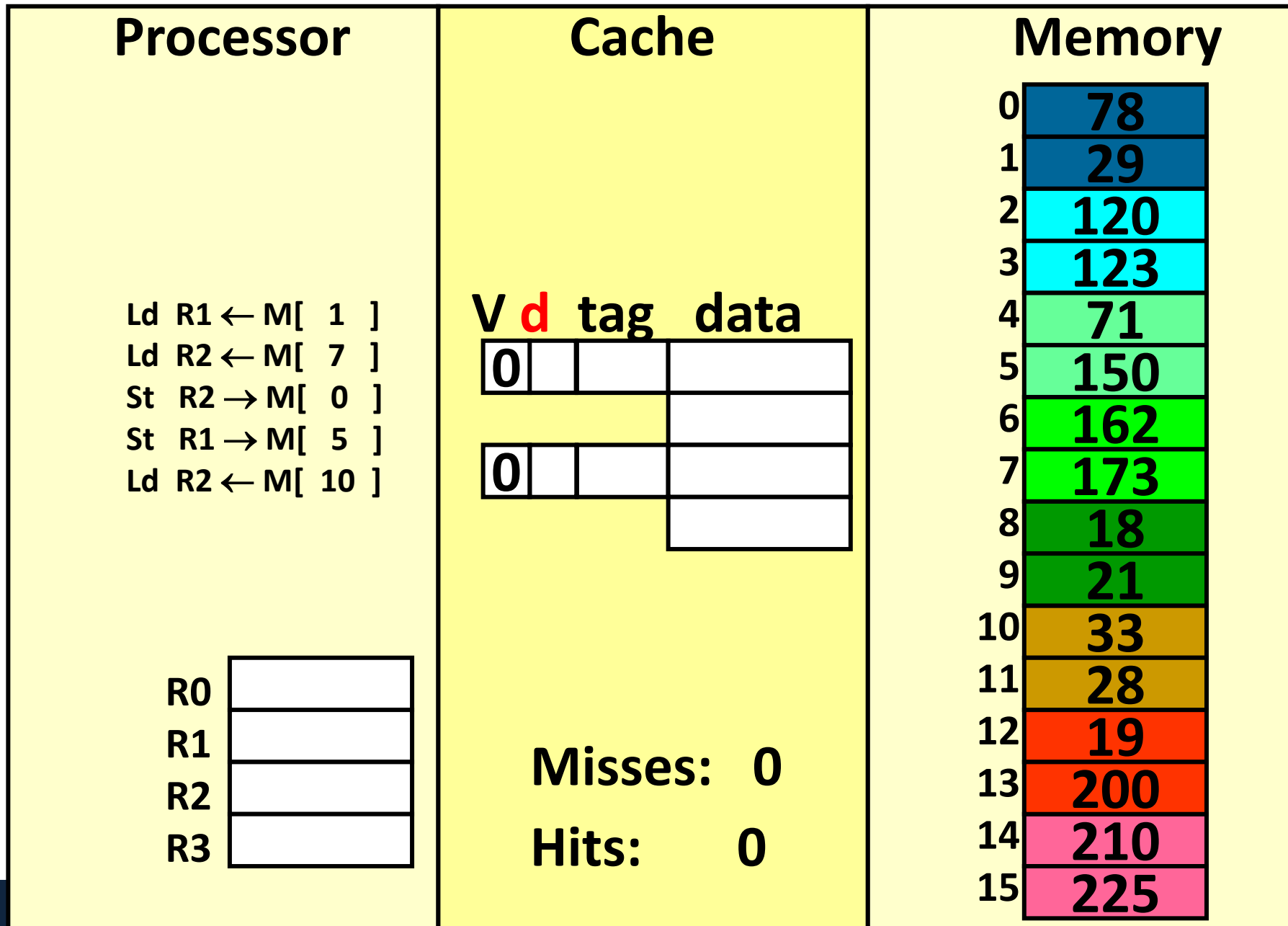
Agenda

- Larger Cache Blocks
- Extra Problems
- LRU with More than Two Blocks
- Write-Through Cache
- **Write-Back Cache**

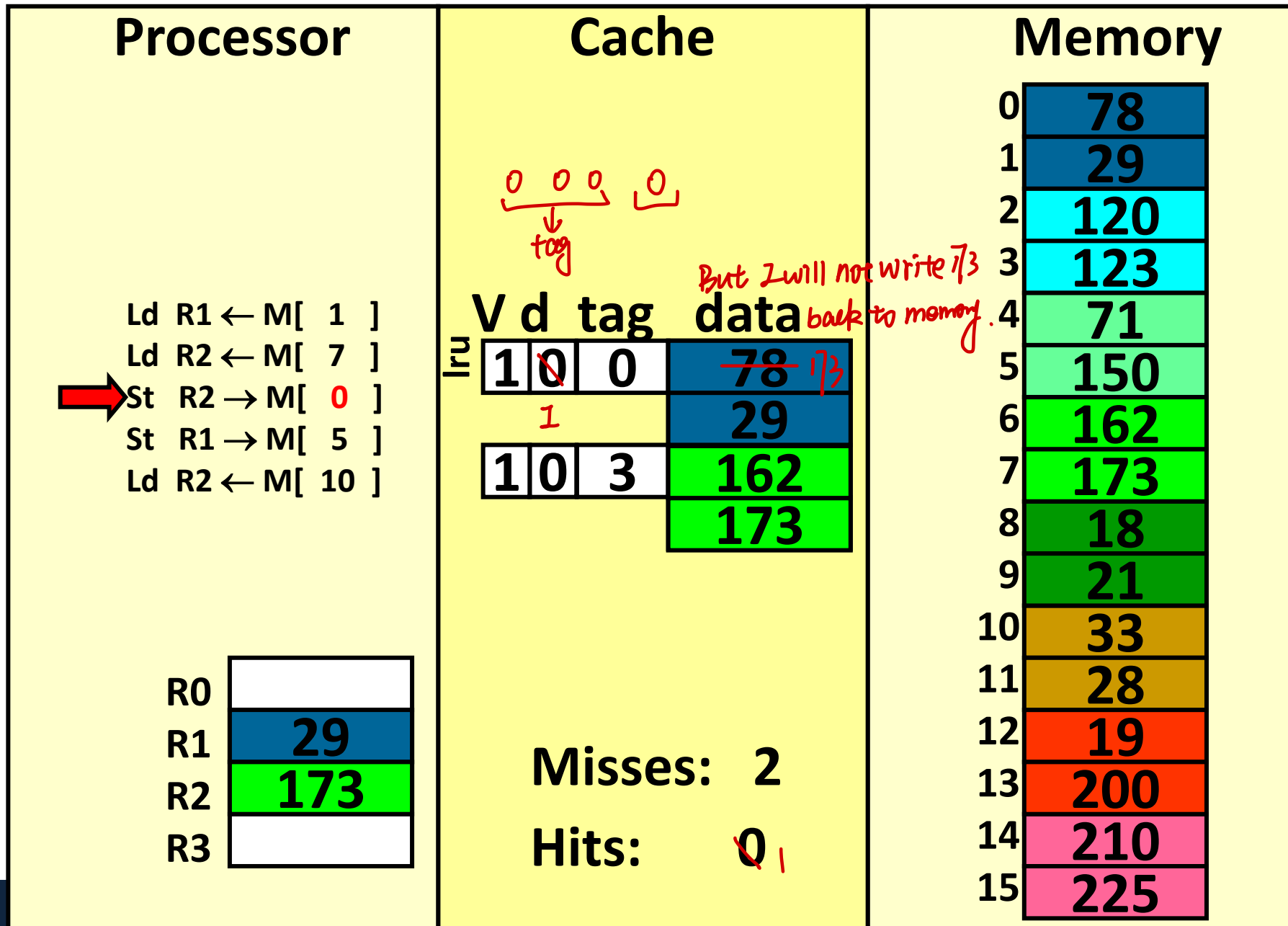
Write-through vs write-back

- Can we design the cache to **NOT** write all stores to memory immediately?
 - Keep the most recent copy in the cache and update the memory **only when** that data is evicted from the cache (**write-back**)
 - Do we need to write-back all evicted lines?
 - No, only blocks that have been modified
 - Keep a “**dirty bit**”, reset when the line is allocated, set when the block is stored into. If a block is “dirty” when evicted, write its data back into memory.

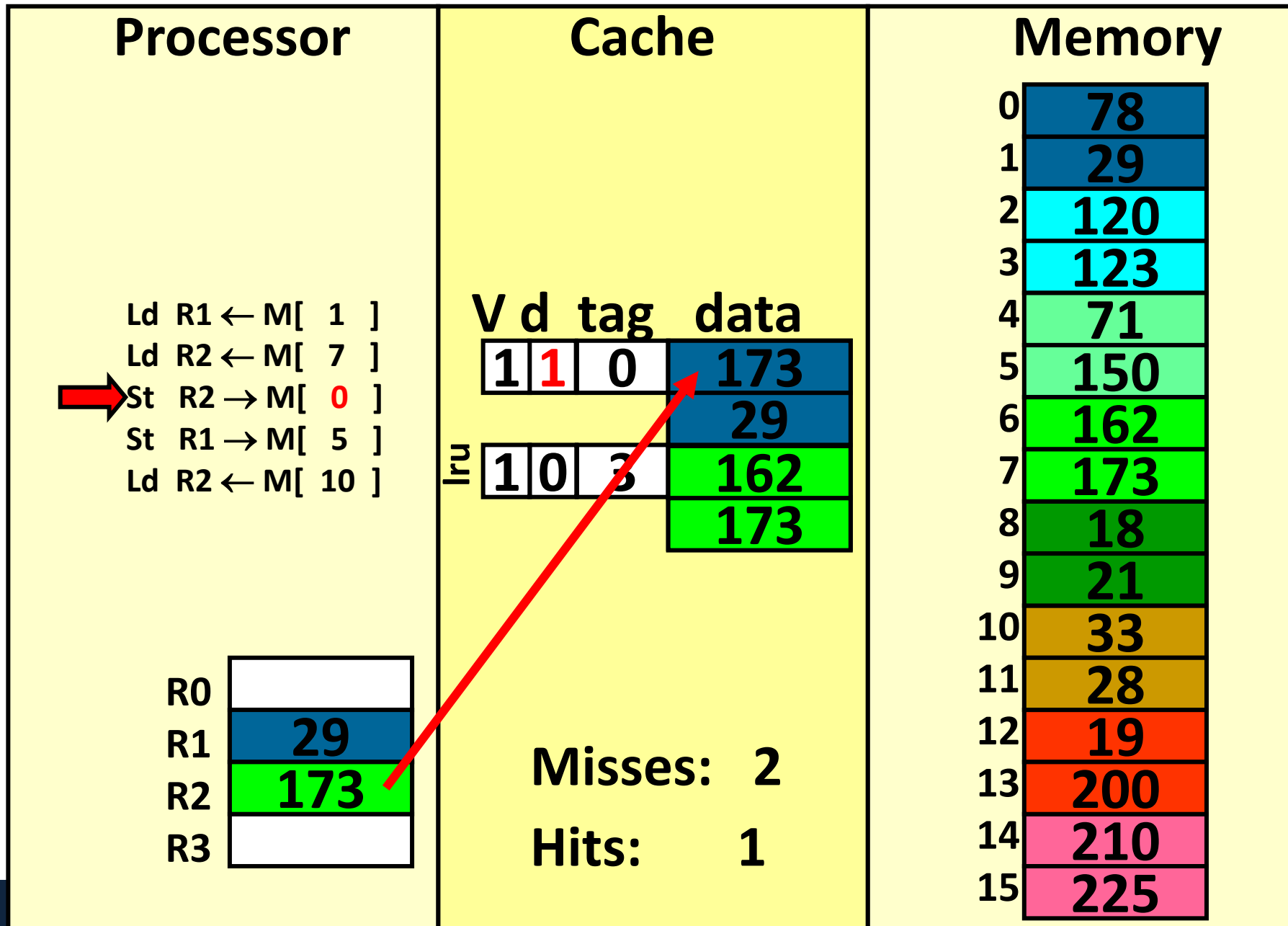
Handling stores (write-back)



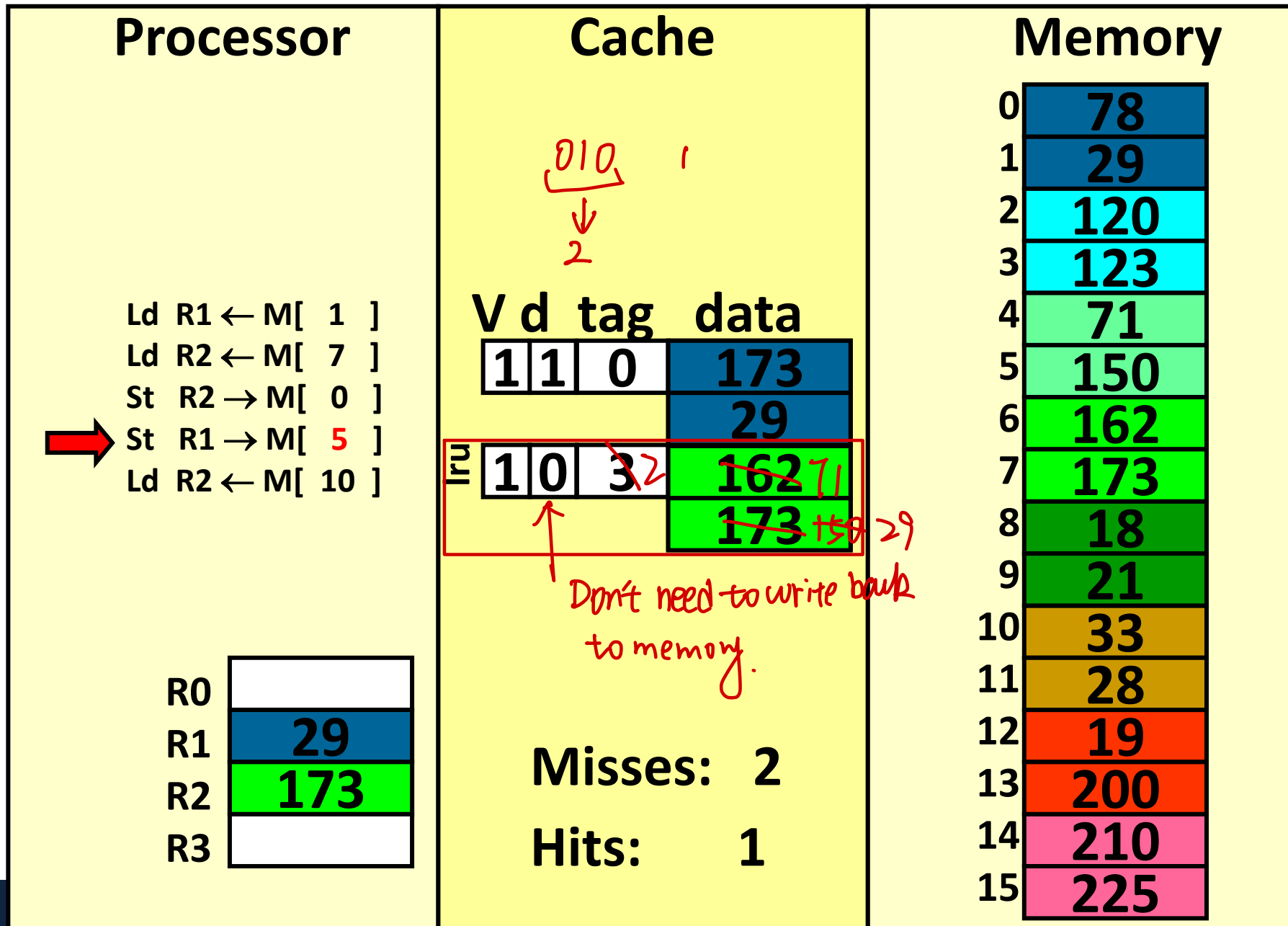
write-back (REF 3)



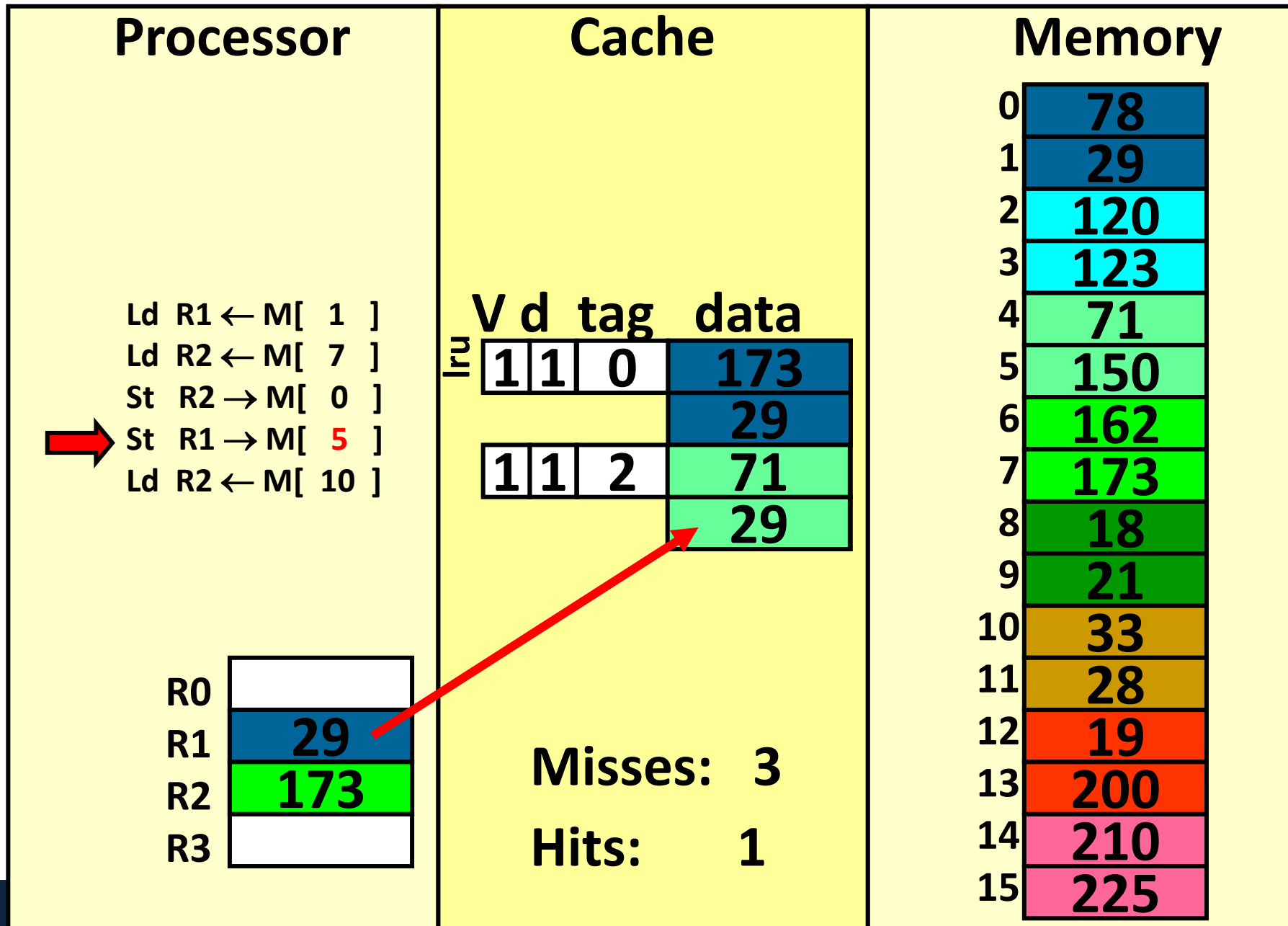
write-back (REF 3)



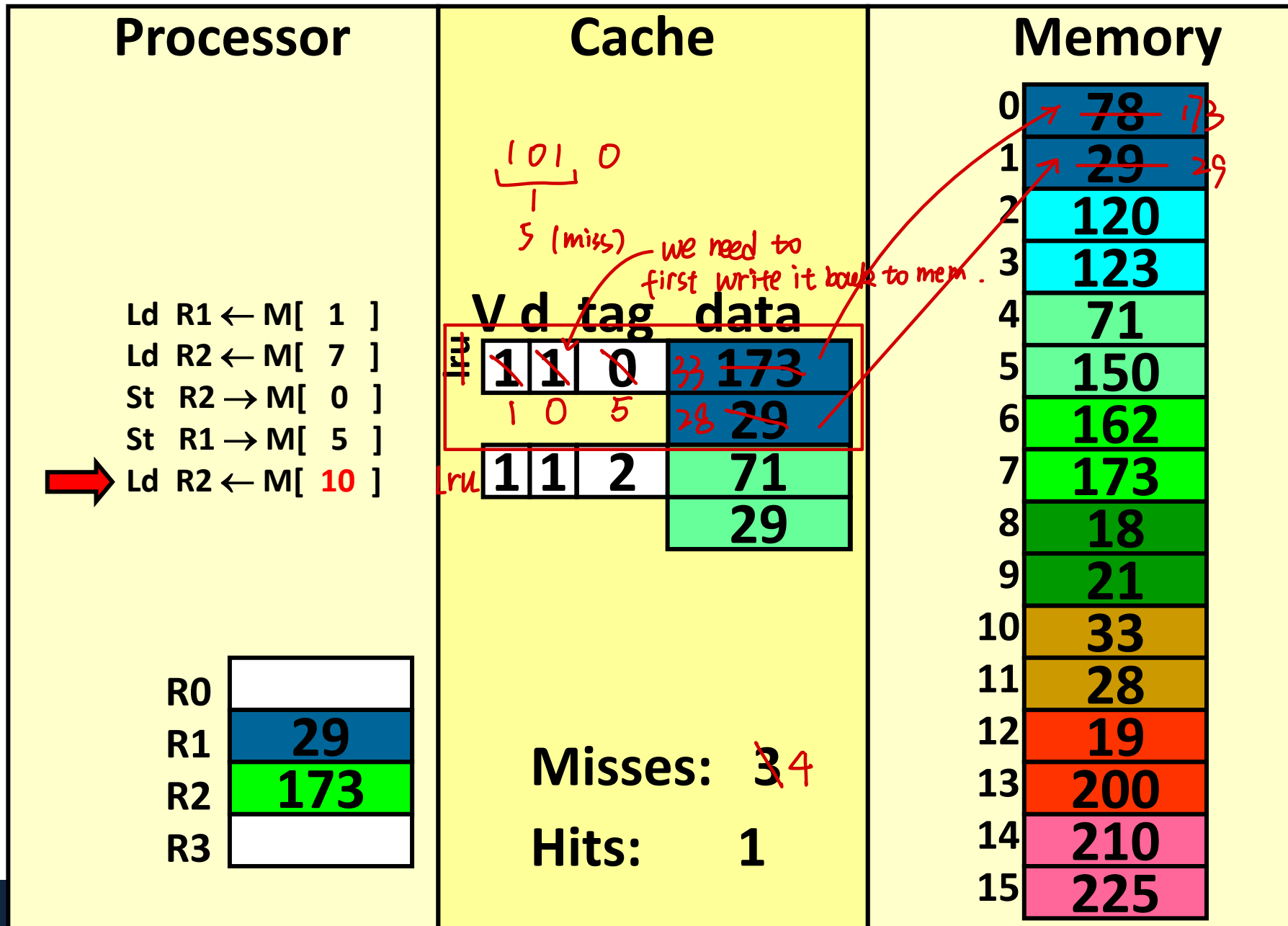
write-back (REF 4)



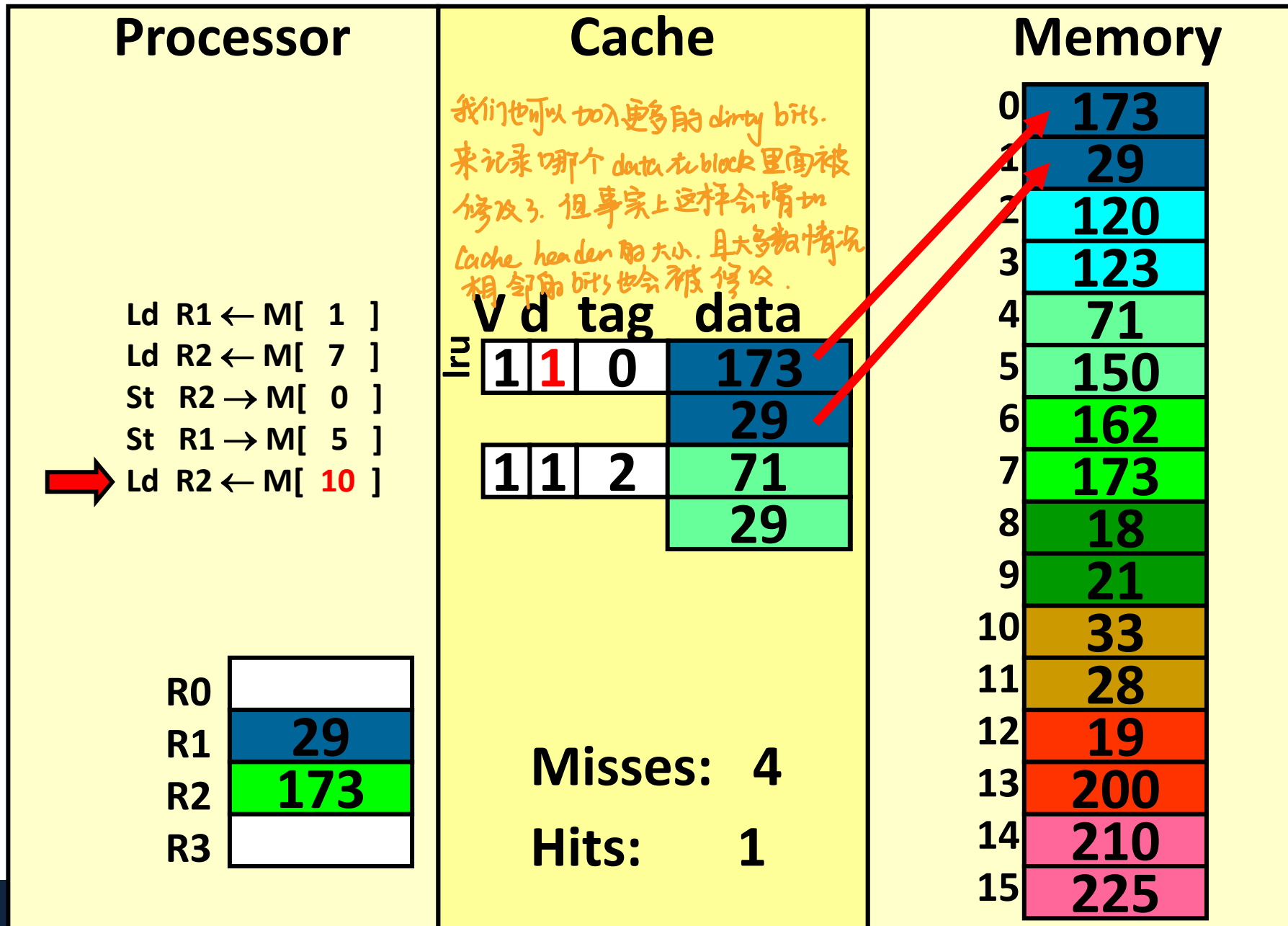
write-back (REF 4)



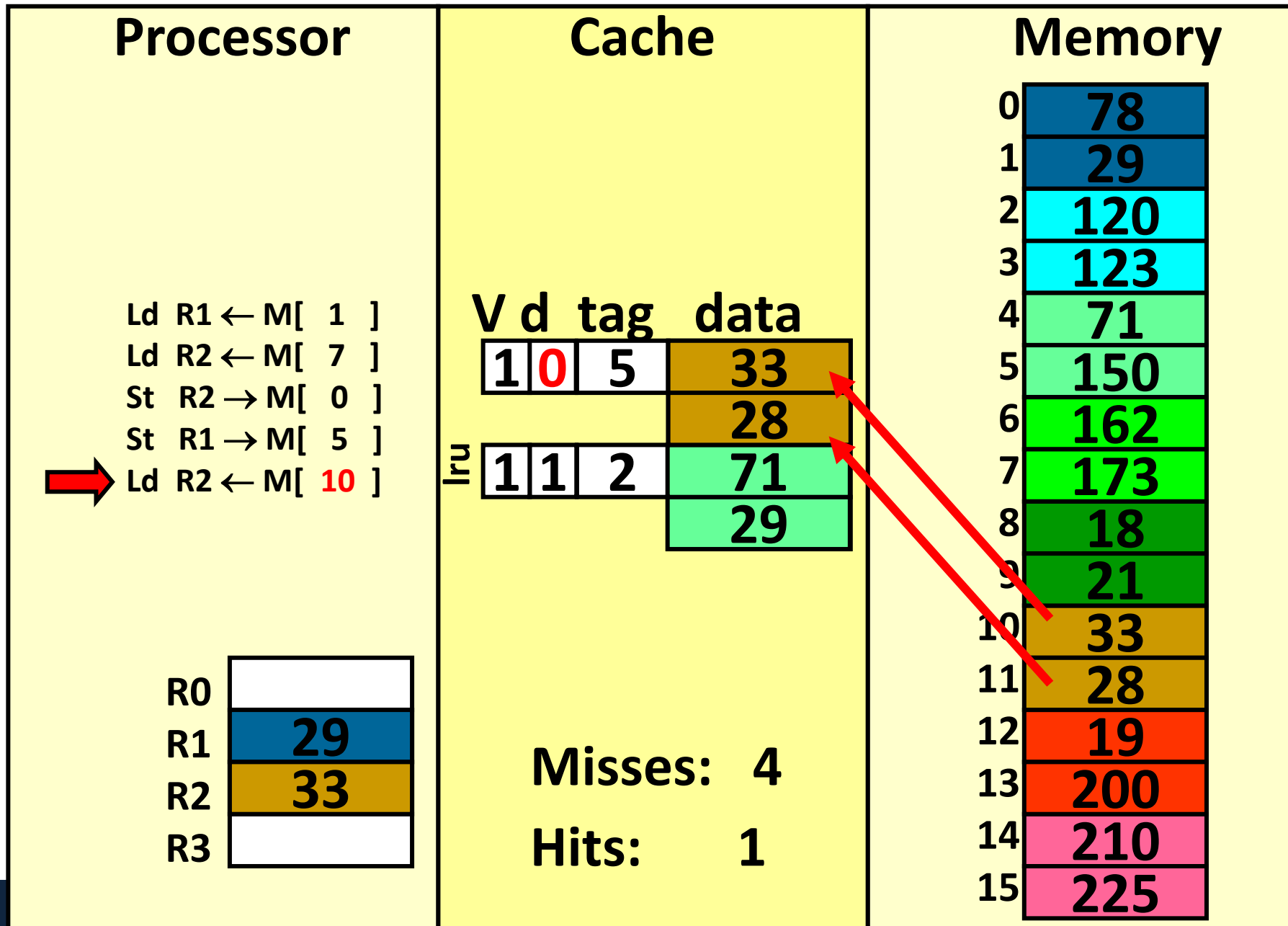
write-back (REF 5)



write-back (REF 5)



write-back (REF 5)



How many memory references?

- Each miss reads a block
 - 2 bytes in this cache
- Each evicted dirty cache line writes a block
- Total reads: 8 bytes
- Total writes: 4 bytes (after final eviction)
*2 dirty block. each has 2 bytes $2*2=4$*

For this example, would you choose write-back or write-through?

Write-back works best when we write to a particular address multiple times before evicting

Review: Writes

Store w No Allocate	Write-Back	Write-Through
Hit?	Write Cache	Write to Cache + Memory
Miss?	Write to Memory	Write to Memory
Replace block?	If evicted block is dirty, write to Memory	Do Nothing
Store w Allocate	Write-Back	Write-Through
Hit?	Write Cache	Write to Cache + Memory
Miss?	Read from Memory to Cache, Allocate to LRU block Write to Cache	Read from Memory to Cache, Allocate to LRU block Write to Cache + Memory
Replace block?	If evicted block is dirty, write to Memory	Do Nothing

Next time

- Direct-mapped vs associative caches.

Next time

- Direct-mapped vs associative caches.